

This Page Is Inserted by IFW Operations  
and is not a part of the Official Record

## **BEST AVAILABLE IMAGES**

Defective images within this document are accurate representations of the original documents submitted by the applicant.

Defects in the images may include (but are not limited to):

- BLACK BORDERS
- TEXT CUT OFF AT TOP, BOTTOM OR SIDES
- FADED TEXT
- ILLEGIBLE TEXT
- SKEWED/SLANTED IMAGES
- COLORED PHOTOS
- BLACK OR VERY BLACK AND WHITE DARK PHOTOS
- GRAY SCALE DOCUMENTS

**IMAGES ARE BEST AVAILABLE COPY.**

**As rescanning documents *will not* correct images,  
please do not report the images to the  
Image Problem Mailbox.**

#3

S&H Form: (2/01)

Attorney Docket No. 826.1730

**IN THE UNITED STATES PATENT AND TRADEMARK OFFICE**

In re Patent Application of:

Hiroshi TSUDA

Application No.:

Group Art Unit:

Filed: June 14, 2001

Examiner:

For: DOCUMENT COLLECTION APPARATUS AND METHOD FOR SPECIFIC USE, AND  
STORAGE MEDIUM STORING PROGRAM USED TO DIRECT COMPUTER TO  
COLLECT DOCUMENTS

J1046 U.S. PTO  
09/880070  
06/14/01

**SUBMISSION OF CERTIFIED COPY OF PRIOR FOREIGN  
APPLICATION IN ACCORDANCE  
WITH THE REQUIREMENTS OF 37 C.F.R. §1.55**

Assistant Commissioner for Patents  
Washington, D.C. 20231

Sir:

In accordance with the provisions of 37 C.F.R. §1.55, the applicant(s) submit(s) herewith  
a certified copy of the following foreign application:

Japanese Patent Application No. 2000-397966

Filed: December 27, 2000

It is respectfully requested that the applicant(s) be given the benefit of the foreign filing  
date(s) as evidenced by the certified papers attached hereto, in accordance with the  
requirements of 35 U.S.C. §119.

Respectfully submitted,

STAAS & HALSEY LLP

Date: June 14, 2001

By: \_\_\_\_\_

James D. Halsey, Jr.  
Registration No. 22,729

700 11th Street, N.W., Ste. 500  
Washington, D.C. 20001

©2001 Staas & Halsey LLP



PATANT OFFICE  
JAPANESE GOVERNMENT

This is to certify that the annexed is a true copy of the  
following application as filed with this Office.

Date of Application: December 27, 2000

Application Number: Patent Application  
No. 2000-397966

Applicant(s): FUJITSU LIMITED

March 9, 2001

Commissioner,  
Patent Office      Kozo OIKAWA

Certificate No. 2001-3017566

日 本 国 特 許 庁  
PATENT OFFICE  
JAPANESE GOVERNMENT

J1046 U.S. PTO  
09/880070  
06/14/01

別紙添付の書類に記載されている事項は下記の出願書類に記載されている事項と同一であることを証明する。

This is to certify that the annexed is a true copy of the following application as filed with this Office.

出 願 年 月 日

Date of Application:

2000年12月27日

出 願 番 号

Application Number:

特願2000-397966

出 願 人

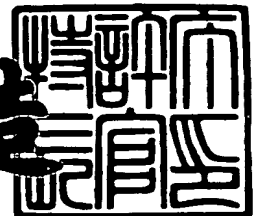
Applicant (s):

富士通株式会社

2001年 3月 9日

特許庁長官  
Commissioner,  
Patent Office

及 川 耕 造



出証番号 出証特2001-3017566

【書類名】 特許願

【整理番号】 0051888

【提出日】 平成12年12月27日

【あて先】 特許庁長官殿

【国際特許分類】 G06F 17/30

【発明の名称】 特定用途向けの文書収集装置、その方法及びコンピュータに実行させるためのプログラムを記録した記録媒体

【請求項の数】 10

【発明者】

【住所又は居所】 神奈川県川崎市中原区上小田中4丁目1番1号 富士通株式会社内

【氏名】 津田 宏

【特許出願人】

【識別番号】 000005223

【氏名又は名称】 富士通株式会社

【代理人】

【識別番号】 100074099

【住所又は居所】 東京都千代田区二番町8番地20 二番町ビル3F

【弁理士】

【氏名又は名称】 大菅 義之

【電話番号】 03-3238-0031

【選任した代理人】

【識別番号】 100067987

【住所又は居所】 神奈川県横浜市鶴見区北寺尾7-25-28-503

【弁理士】

【氏名又は名称】 久木元 彰

【電話番号】 045-573-3683

【手数料の表示】

【予納台帳番号】 012542

【納付金額】 21,000円

【提出物件の目録】

【物件名】 明細書 1

【物件名】 図面 1

【物件名】 要約書 1

【包括委任状番号】 9705047

【プルーフの要否】 要

【書類名】 明細書

【発明の名称】 特定用途向けの文書収集装置、その方法及びコンピュータ  
に実行させるためのプログラムを記録した記録媒体

【特許請求の範囲】

【請求項 1】 ネットワークから文書を収集する文書収集方法であって、  
前記文書の参照関係に基づいて、前記ネットワーク上のコミュニティ内から文  
書を所定数以上収集し、

前記コミュニティ内から前記所定数以上の文書を収集した後、収集済み文書の  
参照関係に基づいて、前記コミュニティ内外から文書を収集する、  
ことを特徴とする文書収集方法。

【請求項 2】 前記収集済み文書の参照関係及びネットワーク上の場所を示  
す情報に基づいて重要度を算出し、

前記参照関係及び前記重要度に基づいて、収集すべき文書を決定する、  
ことを特徴とする請求項 1 記載の文書収集方法。

【請求項 3】 前記収集すべき文書は、前記コミュニティ内外別に決定され  
る、

ことを特徴とする請求項 2 記載の文書収集方法。

【請求項 4】 前記収集済み文書を検索した結果を、前記コミュニティ内外  
に分けて提示する、

ことを特徴とする請求項 3 記載の文書収集方法。

【請求項 5】 ネットワークから文書を収集する文書収集方法であって、  
ある分野に関する文書群である正例文書群と、前記分野と関連が少ない分野  
に関する文書群である負例文書群とを与え、

前記正例文書群及び前記負例文書群の参照関係に基づいて、前記分野に関する  
収集すべき文書を決定し、前記ネットワークから前記収集すべき文書を収集す  
る、

ことを特徴とする文書収集方法。

【請求項 6】 前記参照関係に基づいて、前記正例文書群の文書からのみ参  
照される度合いである参照度を算出し、

前記参照度が高い文書を収集すべき文書として決定する、  
ことを特徴とする請求項 5 記載の文書収集方法。

【請求項 7】 前記参照関係に基づいて、前記正例文書群の文書を参照している収集済み文書から参照されている文書について、収集済み文書からの被参照数を示す共参照度を算出し、

前記共参照度が高い文書を収集すべき文書として決定する、  
ことを特徴とする請求項 5 又は 6 記載の文書収集方法。

【請求項 8】 前記負例文書群は、複数の分野に関する文書群の和集合である、

ことを特徴とする、請求項 5 乃至 7 記載の文書収集方法。

【請求項 9】 前記収集済み文書で用いられている参照表現に基づいて、前記収集済み文書をまとめあげる、

ことを特徴とする請求項 1 乃至 8 記載の文書収集方法。

【請求項 10】 前記収集済み文書で用いられている参照表現に基づいて、前記収集済み文書にキーワードを付与する、

ことを特徴とする請求項 1 乃至 9 記載の文書収集方法。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】

本発明は、文書の収集に関し、特に特定用途に合わせて文書を効率的に収集する文書収集装置、その方法に関する。

【0002】

【従来の技術】

イントラネット、WWW等のネットワーク上の文書の検索エンジンは、ネットワークから文書を収集する文書収集装置（ロボット）と、収集した文書用のキーワード索引を作成する検索エンジンとから実現される。

【0003】

文書収集装置は、所与のネタURL（Uniform Resource Locator）集（収集を開始する際の開始点となるURL集）から文書収集を開始し、収集済みの文書か



らアンカー（参照関係）により参照されている未収集文書を次収集候補として収集し、といった処理を一定の回数繰り返すことにより動作する。このようにして文書収集ロボットは、数千万から数億のURLから文書を定期的に収集する。ここで、URLとは、ネットワーク上の情報のありかと取得方法を指定する記述方式をいう。

## 【 0 0 0 4 】

ところで、今日、ネットワーク上の文書の増加スピードは速く、2000年1月には、Inktomi等によって、インターネットのユニーク文書は10億文書に達したという調査結果が発表されている。また、2000年7月には、アメリカCyveillance社によって、インターネットの大きさは約21億文書であり、2001年にはさらに倍の大きさになると予測されるという調査結果が発表されている。

## 【 0 0 0 5 】

10億URLから文書を収集するとともに、一日100万URLずつ（毎秒約10URL＝40Kバイト）収集したとしても収集し終わるには3年かかることになり、収集し終わった頃には最初の頃に収集した文書の情報は陳腐化してしまう。そこで、用途に合わせて重要度の高い情報だけを効率よく収集する知的文書収集装置が求められていた。

## 【 0 0 0 6 】

特定用途の文書を優先して収集する文書収集装置には、以下のものがある。

- ・例えば、特開平9-311802に開示される発明のように、新しい情報を優先して収集する。
- ・内容が類似していると考えられる文書を収集する。その際に、以下の考え方を導入する。

## 【 0 0 0 7 】

a) 階層数で収集範囲を制限する。

例えば、特開平9-218876に開示される発明のように、参照関係を有する文書は内容的にも近いと考えられるが、あまり階層的に離れると意味的な繋がりがなくなるため、階層数で収集範囲を制限して文書を収集するという考え方。

## 【 0 0 0 8 】

b) 意味的内容が近い文書のみ収集する。

例えば、特開平10-105572 に開示される発明のように、文書の中身のマッチングから意味的な近さを計算し、参照関係を有する文書のうち、意味的に近い文書だけを収集するという考え方。

#### 【0009】

c) 参照先を示す文字列が適当な文書のみ収集する。

例えば、特開平10-260979 及び特開2000-9011 に開示される発明のように、参照先を表している表現である参照表現、例えばHTMLであればアンカータグの内容に基づいて、その参照表現で参照されている参照先文書を次に収集するか否かを判定するという考え方。

- ・一般的に、より人気度の高い文書から優先して収集する。

#### 【0010】

被参照数の多い文書は、人気度が高いと考えられる。収集済みの文書群の文書から参照されている数（被参照数）が多い文書から順に収集することで、重要度の高い文書を優先して収集できるという考え方。

#### 【0011】

##### 【発明が解決しようとする課題】

しかし、上述の従来技術の枠組みだけでは、企業のようなコミュニティのポータルサイトに求められるような文書の収集に用いるためには、不十分な点があった。例えば、企業内のポータルサイト、つまりコーポレートポータルの要件として、以下の点が要求される。

- ・社内外でリアルタイムに発生する膨大な文書を自動的に収集する。
- ・自動で意味解析及び分類分け（カテゴライズ）する。
- ・文書を収集し、分類した結果を画面の適当な場所に（人に合わせて）フィードする。

#### 【0012】

このうち、文書収集において、社内外の膨大な文書を漫然と収集するのではなく、文書の中から業務に関係するという観点から文書を選別して収集することが必要とされる。業務に関係するという観点は、特定の意味的内容を持つ、或いは

重要度を持つということとはやや異なる。例えば、ある程度の規模の企業が有するイントラネットコミュニティでは、文書内容も意味的に多様になるからである。また、社外（例えばインターネット）の文書は、趣味に関する情報も人気度が高くそうした情報は必ずしもコーポレートポータルにとって有用であるとは限らない。

#### 【0013】

しかし、従来の文書収集において用いられてきた枠組み、例えば、最新情報の優先取得、特定分野情報の優先取得、重要度優先取得という枠組みだけでは、このような趣味に関する情報のように、一般的に重要度が高いが必ずしもこのコミュニティにとって有用でない文書も収集されてしまうという問題があった。

#### 【0014】

また、例えば、上述の従来技術の「意味的内容が近い文書のみを収集する」と方法で文書を収集する場合、各々の考え方には以下の問題があった。

- ・単に階層数を予め制限する考え方は、処理は簡単であるが、本当に意味内容が近い文書を優先して収集しているのか、また、重要な文書を収集し逃していないのか、保証がない。

- ・文書の内容を比べて意味的内容が近いかな否か判定する方式によれば、一般に自然言語処理を使って、文書に記載された本文を解析してキーワードを取り出し、取り出されたキーワードの類似度によって解析する。そのため、処理に時間がかかる。早くても、毎秒100文書程度しか処理できない。従って、数十億ともいわれる文書を1つ1つ処理することは、現実的な時間内に行いがたい。また、そのように時間をかけて処理したとしても、その精度は70から80%程度である。さらに、この処理は、言語の種類に大きく依存するため、言語毎に判定ツールを備えることが必要となる。

- ・参照表現に基づいて収集するか否か判定する場合でも、参照表現で用いられる文字列には、「ホームページ」、「トップに戻る」及び「ここをクリック」といったような決まった語句（定番的ば語句）も多く、必ずしも参照先の意味的内容を表しているとは限らない。

#### 【0015】

以上の問題を鑑み、用途にあった文書を言語に依存せず、かつ精度良く迅速に収集することを可能とすることが、本発明が解決しようとする課題である。

【0016】

【課題を解決するための手段】

本発明は、ネットワークから文書の収集を行なう装置または方法を前提とする。そして、本発明の各態様に係る装置では、ネットワークから文書を収集する文書収集装置において、収集済みの文書群の参照関係に基づいて、次に収集すべき文書の候補である次収集候補を決定する次候補判定手段と、ネットワークから前記次収集候補を収集して収集済み文書群に加える文書収集手段と、を備え、収集済み文書群の文書がある数以上になるまで、次候補判定手段による次収集候補の決定及び前記文書収集手段による文書の収集を繰り返す。

【0017】

上記装置を、ネットワーク上のコミュニティにとって有用度の高い文書を収集するコミュニティ向けの文書収集装置として構成するようにしてもよい。そのために、上記構成において、文書収集手段がネットワーク上のコミュニティ内から文書をまんべんなく収集した後、次候補判定手段は、収集済み文書群の参照関係に基づいてコミュニティ内外の文書から次収集候補を決定する、こととしてもよい。

【0018】

コミュニティ内外から文書を収集する前に、コミュニティ内から文書をまんべんなく収集することにより、コミュニティ内で必要とされている多様な分野の文書についての情報を入手することができる。このようにして入手した多様な分野に関する文書群の参照関係を用いてコミュニティ内外から文書を収集することにより、正確にコミュニティにとって有用度の高い文書を収集することが可能となる。また、文書本文の内容を解析しないため、言語に依存せず、迅速にコミュニティにとって有用度の高い文書を収集することが可能となる。

【0019】

上記構成において、収集済み文書群の参照関係及び文書のネットワーク上の場所を示す情報、例えばURL、に基づいて重要度を算出するランキング手段を更

に備え、次候補判定手段は、参照関係及び重要度に基づいて次収集候補を決定することとしてもよい。

【 0 0 2 0 】

上記コミュニティ向け文書収集装置において、ランキング手段は、重要度に基づいて、前記コミュニティ内外に分けてランキングし、次候補判定手段は、コミュニティ内及びコミュニティ外それぞれにおいて、ランキングが高い文書を前記次収集候補とすることとしてもよい。これにより、次収集候補がコミュニティ内又はコミュニティ外に集中し、文書がコミュニティ内又はコミュニティ外いずれかからばかり収集されてしまうことを防ぐことが可能となる。

【 0 0 2 1 】

また、上記コミュニティ向け文書収集装置は、更に、収集済み文書群を検索した結果を、前記コミュニティ内外に分けて提示する提示手段を備えることとしても良い。これにより、コミュニティに属するクライアントが、コミュニティ内外別に文書の検索結果を取得することが可能となる。

【 0 0 2 2 】

また、上記コミュニティ向け文書収集装置は、更に、文書がコミュニティ内の文書であるか否かを文書のネットワーク上での場所を示す情報、例えばURL、に基づいて判別するコミュニティ判別手段を備えることとしても良い。文書のネットワーク上での場所を示す情報に基づいて判定することにより、文書がコミュニティ内の文書であるか否かの判定が迅速に行うことが可能となる。

【 0 0 2 3 】

また、上記のネットワークから文書を収集する文書収集装置を、特定の分野に関する文書を収集する特定分野向け文書収集装置として構成するようにしてもよい。そのために、本発明の更なる別の態様によれば、ネットワークから文書を収集する装置において、文書の収集に先立って、特定分野に関する文書群である正例文書群と、特定分野と関連が少ない分野に関する文書群である負例文書群とを収集済み文書群として与え、文書収集手段は、収集された次収集候補を、正例文書群に加え、収集済み文書群のうち、正例文書群の文書がある数以上になるまで、次候補判定手段による次収集候補の決定及び文書収集手段による収集を繰り返

すように構成する。これにより、特定分野に関する文書を、文書本文の内容を解析せずに、参照関係に基づいて迅速に収集することが可能となる。

【 0 0 2 4 】

また、上記の特定分野向け文書収集装置において、更に、収集済み文書の参照関係に基づいて、正例文書群の文書からのみ参照される度合いである参照度を算出する参照度算出手段を備え、次候補判定手段は、参照度が高い文書を次収集候補として決定することとしてもよい。また、上記の特定分野向け文書収集装置において、更に、収集済み文書の参照関係に基づいて、正例文書群の文書を参照している収集済み文書群から参照されている文書について、収集済み文書群からの

被参照数を示す共参照度を算出する共参照度算出手段を備え、次候補判定手段は、共参照度が高い文書を次収集候補として決定することとしてもよい。参照度及び共参照度を用いることにより、収集したい分野に関する文書を、文書本文の内容を検討すること無く、迅速に収集することが可能となる。

【 0 0 2 5 】

また、上記の特定分野向け文書収集装置は、複数の分野を対象とし、各分野に関する文書を同時に収集する文書収集装置とすることもできる。そのために、上記の特定分野向け文書収集装置において、収集に先立って与える収集済み文書群を複数の分野に関する文書群の和集合とし、ある分野に関する文書群を正例文書群として文書を収集する際に、他の残りの分野に関する文書群の和集合を負例文書群とするように構成する。

【 0 0 2 6 】

また、各文書収集装置において、更に、収集済み文書で用いられている参照表現に基づいて収集済み文書群をまとめあげるまとめあげる手段を更に備えることとしてもよい。参照表現のうち、参照先文書と参照元文書の内容が同一であるのにネットワーク上で分散されて格納されていることを示す参照表現がある。例えば、「次へ」、「Next」、「前へ」及び「Prev」等がそのような参照表現に該当する。まとめあげ手段は、このような参照表現による参照関係をもつ2つ以上の文書を1つにまとめあげる。

## 【 0 0 2 7 】

また、各文書収集装置において、更に、収集済み文書群の文書である収集済み文書で用いられている参照表現に基づいて、収集済み文書にキーワードを付与するキーワード付与手段を備えることとしても良い。これにより、文書本文の意味内容を解析することなく、かつ、様々な各キーワードの異称をも、キーワードとすることが可能となる。

## 【 0 0 2 8 】

また、キーワード付与手段は、参照表現が参照先文書に関係なく使用される参照表現の場合、キーワードとしないこととしても良い。ここで、参照先文書に関係なく使用される参照表現の例として、「トップへ戻る」、「ホームへ」等が考えられる。

## 【 0 0 2 9 】

また、キーワード付与手段は、参照表現が参照する相異なる文書数を計数し、相異なる文書数がある数以上である場合、その参照表現をキーワードとしないこととしても良い。このような参照表現は、参照先文書に関係なく使用される参照表現である可能性が高いからである。

## 【 0 0 3 0 】

また、キーワード付与手段は、参照表現が参照する相異なる文書数がある数未満である場合、更に、各収集済み文書でその参照表現が参照されている回数である参照回数を計数し、相異なる文書数及び参照回数に基づいて、その参照表現をキーワードとするか否か判定することとしてもよい。

## 【 0 0 3 1 】

また、キーワード付与手段は、参照表現に基づくキーワードに、収集済み文書の本文から抽出したキーワード及び収集済み文書のURLから抽出したキーワードを組み合わせることとしてもよい。これにより、多様な方法で抽出したキーワードを組み合わせることが可能となる。

## 【 0 0 3 2 】

また、本発明の各構成により行われる処理の過程からなる方法によっても、前述した課題を解決することができる。また、上述した本発明の各構成により行な

われる機能と同様の制御をコンピュータに行なわせるプログラムを記憶したコンピュータ読み取り可能な記憶媒体から、そのプログラムをコンピュータに読み出させて実行させることによっても、前述した課題を解決することができる。

#### 【 0 0 3 3 】

##### 【発明の実施の形態】

以下、本発明の実施の形態を図面に基づいて説明する。本発明は、ネットワークから、用途にあった文書を収集する文書収集装置に関する。図 1 に、本発明の原理図を示す。図 1 に示すように、文書収集装置 1 は、インターネットやイントラネット等のネットワークに接続されている。文書収集装置 1 は、文書収集手段 2、参照関係抽出手段 3、コミュニティ判別手段 4、次候補判定手段 5、ランキング手段 6、URL 判定手段 7、参照度／共参照度算出手段 8、まとめあげ手段 9、キーワード抽出手段 10 を備える。図 1 において、点線で示される構成要素、つまり、コミュニティ判別手段 4 及び参照度／共参照度算出手段 8 は、実施形態によって用いられったり、用いられなかったりする。同様に、点線で示される矢印、つまり、ランキング手段 6 による文書のランキング結果は、実施形態によって、次候補判定手段 15 による次収集候補の判定に用いられったり、用いられなかったりする。

#### 【 0 0 3 4 】

本発明の 1 実施形態に係わる文書収集装置は、ネットワーク上のコミュニティ向けの文書を収集する。そのために、1 実施形態に係わるコミュニティ向け文書収集装置は、文書収集手段 2、参照関係抽出手段 3、コミュニティ判別手段 4、次候補判定手段 5、ランキング手段 6、まとめあげ手段 9 及びキーワード付与手段 10 を備える。コミュニティ向け文書収集装置において、まず、コミュニティ内からまんべんなく文書を収集した後、コミュニティ内外からコミュニティにとって有用度が高い文書を収集する。

#### 【 0 0 3 5 】

参照関係抽出手段 3 は、収集済み文書群 20 から参照関係を抽出し、文書間参照関係 22 を抽出する。なお、収集開始時は、予め収集済み文書群 20 として初期文書群を与える。コミュニティ判別手段 4 は、収集済み文書群 20 の参照先文



書であって、未収集の文書がコミュニティ内の文書であるか否か判別する。

【 0 0 3 6 】

次候補判定手段 5 は、収集済み文書群 2 0 の参照先であって、コミュニティ内の未収集文書を次収集候補 2 1 として判定する。文書収集手段 2 は、次収集候補 2 1 として判定された文書を収集し、新たに収集した文書群（新規収集文書群）を収集済み文書群 2 0 に加え、新たな収集済み文書群 2 0 とする。文書収集手段 2 は、収集済み文書群 2 0 の文書数が規定された値以上であるか否か判定する。収集済み文書群 2 0 の文書数が規定された値より少ない場合、上述のようにしてコミュニティ内から文書を収集する処理を繰り返す。このようにコミュニティ内の文書を規定数以上、まんべんなく収集することにより、コミュニティ内の文書が属する多様な分野についての情報を取得する。この情報は、コミュニティにとって有用度が高い文書をコミュニティ内外から収集することに役立てられる。

【 0 0 3 7 】

収集済み文書群 2 0 の文書数が規定された値以上である場合、次にコミュニティにとって有用度が高い文書をコミュニティ内外から収集する。参照関係抽出手段 3 により新規収集文書群から参照関係を抽出し、コミュニティ判別手段 4 により参照先文書であって未収集の文書がコミュニティ内の文書であるか否か判別する。ランキング手段 6 は、参照関係及び、文書のネットワーク上での場所を示す情報、例えば URL、の特徴に基づいて、収集済み文書の参照先となっている未収集の文書をコミュニティ内外別にランキングする。ランキング手段 6 は、URL 判定手段 7 を備え、URL 判定手段 7 は、参照先文書と参照元文書の URL 文字列上の類似を判定する。ランキング手段 6 は、URL 判定手段 7 によって判定され URL の文字列上の類似に基づいて、文書をランキングする。

【 0 0 3 8 】

次候補判定手段 5 は、コミュニティ内外でそれぞれ上位にランキングされた未収集文書を次回にネットワークから収集すべき文書である次収集候補 2 1 として判定し、文書収集手段 2 は、次収集候補 2 1 として判定された文書を収集する。このように、本発明の 1 実施形態に係わるコミュニティ向け文書収集装置は、多段階に分けてコミュニティにとって有用度が高い文書を収集する。ある規定され

た以上の文書をコミュニティ内外から収集すると、まとめあげ手段 9 は、参照表現に基づいて収集済み文書 2 0 をまとめあげる。キーワード付与手段 1 0 は、参照表現及び参照表現の出現頻度に基づいて、収集済み文書 2 0 にキーワードを付与する。ランキング手段 6 は、上述のようにして、今度は収集済み文書 2 0 をランキングする。最終的にまとめあげられ、キーワードを付与し、ランキングした収集済み文書 2 0 は、収集文書ファイル 2 3 として格納される。上述のように、コミュニティ向け文書収集装置において、文書本文の内容を解析していないため、言語に依存せず、迅速に、用途に合った文書を収集することができる。

## 【 0 0 3 9 】

また、本発明の別の 1 実施形態に係わる文書収集装置は、特定の分野に関する文書を収集する。そのために、上記特定分野に関する文書収集装置は、文書収集手段 2、参照関係抽出手段 3、次候補判定手段 5、ランキング手段 6、参照度／共参照度算出手段 8、まとめあげ手段 9 及びキーワード付与手段 1 0 を備える。特定分野に関する文書収集装置において、コミュニティ内外の文書の区別は不要であるため、コミュニティ判別に係わる処理はない。

## 【 0 0 4 0 】

特定分野に関する文書収集装置において、収集に先立って特定分野に関する文書群を正例文書群として、その特定分野との関連が少ない文書群を負例文書群として与える。収集済み文書群 2 0 は、正例文書群と負例文書群の和集合とする。参照度／共参照度算出手段 8 は、ある文書と正例文書群、その文書と負例文書群のそれぞれの参照関係に基づいて、その文書が特定分野に関連する度合いを参照度及び共参照度として算出する。次候補判定手段 5 は、ランキング手段 6 によるランキングの代わりに、参照度／共参照度算出手段 8 が算出した参照度又は共参照度が高い未収集文書を次収集候補として判定する。また、負例文書群に含まれる収集済み文書 2 0 のうち、参照度又は共参照度が高い文書を負例文書群から除き、正例文書群に加える。文書収集手段 2 は、次収集候補 2 1 として判定された文書を収集し、正例文書群に加える。そして、正例文書群の文書数が規定された数以上になるまで、次収集候補の決定及び文書の収集を繰り返す。その他の動作は、上述の通りである。

## 【 0 0 4 1 】

以下、第 1 実施形態に係わる、コミュニティにとって有用度の高い文書を収集するコミュニティ向け文書収集装置について説明する。本発明の第 1 実施形態において述べるネットワーク上のコミュニティとして、例えば、社内サイト、業界サイト及び特定トピックのネットワーク上のユーザグループが考えられる。ここで、社内サイトは、しばしばイントラネットに代表される。業界サイトは、複数の会社のシステムからなるエクストラネットに代表される。なお、社内サイトに必要な文書を収集する文書収集装置は、コーポレートポータル（E I P : Enterprise Information Portal ともいわれる）ともいわれる、企業内のイントラネットポータルに適用可能である。

## 【 0 0 4 2 】

コミュニティのポータルにおいて、コミュニティにとって有用度が高い文書を優先して自動収集するという要件が必要とされている。例えば、コーポレートポータルの場合、業務に関係する文書を自動収集する必要がある。本発明の第 1 実施形態によれば、このような文書の自動収集を実現する。そのために、第 1 実施形態に係わる文書収集装置において、以下の考え方を採用する。

・特定のコミュニティにとって有用度の高い文書は、そのコミュニティの文書の多くからよく参照されている文書である、またはコミュニティの重要文書から参照されている文書である、と考える。

## 【 0 0 4 3 】

図 2 は、第 1 実施形態に係わる文書収集装置の構成を示す。図 1 に示すように、文書収集装置 1 0 0 は、文書収集部 1 0 1、参照関係抽出部 1 0 2、コミュニティ判別部 1 0 3、次候補判定部 1 0 4、ランキング部 1 0 5、まとめあげ部 1 0 6 及びキーワード付与部 1 0 7 を備える。

## 【 0 0 4 4 】

上述のように、本文書収集装置 1 0 0 において、先にコミュニティ内の文書について複数回、収集を行い、次に、コミュニティ内外の文書についても複数回、収集を行う。このように多段階に分けて複数回、文書収集を行うことが本文書収集装置 1 0 0 の特徴の 1 つである。

## 【0045】

収集開始に先立って、まず、初期文書群を収集済み文書群Sとして与える。この初期文書群は、収集の開始点となる。初期文書群として、例えば、サイトのトップ文書やトップ文書の参照集等が考えられる。収集済み文書群S又は初期文書群は、具体的には、URLテーブル120として文書収集装置100に備えられる。

## 【0046】

続いて、参照関係抽出部102は、収集済み文書群Sから参照関係を抽出し、収集済み文書群Sの参照先となる文書（以下、参照先文書という）のURLをURLテーブル120に格納し、抽出された参照関係を参照関係テーブル121に格納する。コミュニティ判別部103は、参照関係抽出部102が抽出した、収集済み文書群Sの参照先文書が、コミュニティ内の文書であるのか、コミュニティ外の文書であるのか、URLに基づいて判定し、判別結果を参照関係テーブル121に格納する。

## 【0047】

本文書収集装置100は、先にコミュニティ内の文書について1回以上収集を行う。この際、収集をまんべんなく行う。次候補判定部104は、参照関係抽出部102が抽出した収集済み文書群Sの参照先文書のうち、まだ収集されていない、コミュニティ内の文書を次に収集すべき文書の候補（以下、次収集候補Nという）として判定する。文書収集部101は、次収集候補Nとして判定された文書群を収集し、収集した文書を収集済み文書群に追加し、新たな収集済み文書群Sとする。このコミュニティ内の文書の収集は、規定された数の文書を収集するまで行う。コミュニティ内の全ての文書を収集しなくても良く、大体、コミュニティ内の全文書の1/2から1/4程度で良い。まんべんなくコミュニティ内の文書を収集することにより、コミュニティ内で有用な文書の分野についての情報を入手する。

## 【0048】

文書収集部101がコミュニティ内の文書を規定された数だけ収集した後、文書収集装置100は、次に、コミュニティ内外の文書についても1回以上収集を

行う。この場合、上述のようにして、文書収集部 1 0 1 は、文書を収集し、参照関係抽出部 1 0 2 及びコミュニティ判別部 1 0 3 は、URL テーブル 1 2 0 及び参照関係テーブル 1 2 1 に情報を格納した後、さらに、ランキング部 1 0 5 は、参照関係及び文書の URL に基づいて、参照先文書に重要度を与え、その重要度に基づいて、参照先文書をランキングする。

#### 【 0 0 4 9 】

候補判定部 1 0 4 は、ランキング部 1 0 5 による判定結果に基づいて、まだ収集されていない参照先文書であって、コミュニティ内の文書のうちで上位  $n$  1 位内にある文書群、及び、コミュニティ外の文書のうちで上位  $n$  2 位内にある文書群を次収集候補  $N$  となる文書として判定する。コミュニティ内外で分けて次収集候補  $N$  を決定することにより、コミュニティ内とコミュニティ外のいずれかに文書が偏って収集されてしまうことを防ぐことが可能となる。

#### 【 0 0 5 0 】

続いて、コミュニティ内の文書の収集と同様にして、文書収集部 1 0 1 は、次収集候補  $N$  をコミュニティ内外から収集し、収集した文書を収集済み文書群に追加して新たな収集済み文書群  $S$  とする。文書収集装置 1 0 0 は、規定された数の文書を収集するまで、コミュニティ内外からの文書収集を繰り返す。

#### 【 0 0 5 1 】

文書収集部 1 0 1 がコミュニティ内外から規定数だけの文書を収集した後、収集した文書の選別を行う。文書の選別は、まとめあげ部 1 0 6、キーワード付与部 1 0 7 及びランキング部 1 0 5 により行われる。まず、まとめあげ部 1 0 6 は、文書において他文書を参照する際に用いる文字列（参照表現ともいう）に基づいて、収集済み文書のうち、同一内容であるが複数の文書に分割されてい文書をまとめあげる。

#### 【 0 0 5 2 】

キーワード付与部 1 0 7 は、文書中の参照表現に基づいて、キーワードを決定し、文書にキーワードを付与する。より具体的には、キーワード付与部 1 0 7 は、参照表現のうち、「トップに戻る」、「ホームへ」というような参照先に関係なくしばしば使用される参照表現を除く。続いて、キーワード付与部 1 0 7 は、

各参照表現が参照する相異なる文書数を計数し、参照表現テーブル 1 2 2 に格納する（図 2 では不図示）。また、各参照表現について各収集済み文書での出現頻度を計数し、参照回数テーブル 1 2 3 に格納する（図 2 では不図示）。キーワード付与部 1 0 7 は、これら計数結果に基づいて各収集済み文書について参照表現の重みを算出し、重みが大きい順にある数だけの参照表現をキーワードとして各収集済み文書に付与する。

#### 【 0 0 5 3 】

ランキング部 1 0 5 は、参照関係及び文書の URL に基づいて、各文書に重要度を付与し、その重要度に基づいて文書をランキングする。このように、本実施形態に係わるコミュニティ向け文書収集装置 1 0 0 は、文書本文の内容を解析すること無く、参照関係及び URL に基づいて文書を収集し、まとめあげ、キーワードを付与し、ランキングする。

#### 【 0 0 5 4 】

上述のようにして、文書収集装置 1 0 0 は、まとめあげられ、キーワードが付与され、ランキングされた文書群を優良コンテンツ 1 3 0 として提供する。優良コンテンツ 1 3 0 は、検索エンジン 1 4 0 を介して索引 1 4 1 として提供されたり、検索エンジン 1 4 0 を介してサーバ 1 6 0 に提供されたり、分類エンジン 1 5 0 によってディレクトリ編集されてサーバ 1 6 0 に提供されたりする。サーバ 1 6 0 のクライアントは、サーバ 1 6 0 に提供された優良コンテンツ 1 3 0 を、ブラウザ 1 7 0 を介して閲覧することができる。

#### 【 0 0 5 5 】

以下、図 3 から図 6 を用いて各テーブルのデータ構造について説明する。図 3 に URL テーブル 1 2 0 のデータ構造の一例を示す。図 3 に示すように、URL テーブルは、各文書について文書を識別する文書 ID (Identification information)、文書の URL、収集済みであるか否かを示す収集済みフラグ、コミュニティ内の文書であるか否かを示すコミュニティフラグ及び文書の重要度を格納する。文書 ID 及び URL は、参照関係抽出部 1 0 2 が収集済み文書の参照先文書を抽出した際に格納される。収集済みフラグは、文書収集部 1 0 1 がその文書を収集した際に「オン (1)」にされる。コミュニティフラグは、コミュニティ判

別部 103 がその文書がコミュニティ内の文書であると判定した場合に「オン (1)」にされる。重要度は、ランキング部 105 が文書の参照関係及び URL の文字列上の特徴に基づいて算出し、格納する。

## 【0056】

図 4 に参照関係テーブル 121 のデータ構造の一例を示す。図 4 に示すように、参照関係テーブル 121 は、文書の参照関係に関する情報を格納する。より具体的には、参照関係テーブル 121 は、参照元文書の文書 ID である参照元文書 ID、参照元文書によって参照されるコミュニティ内の文書の文書 ID である参照先文書 ID<sub>1</sub>、及び、参照元文書によって参照されるコミュニティ外の文書の文書 ID である参照先文書 ID<sub>2</sub> を格納する。これら情報は、参照関係抽出部 102 によって格納される。

## 【0057】

図 5 に参照表現テーブル 122 のデータ構造の一例を示す。図 5 に示すように、参照表現テーブル 122 は、収集済み文書で各参照表現が用いられる頻度に関する情報を格納する。より具体的には、参照表現テーブル 122 は、各参照表現について、参照表現を識別する表現 ID、参照表現 (文字列)、参照表現が参照する相異なる文書の数である頻度 DF (w)、及び、キーワードとして用いるべきか否かを示す要否フラグを格納する。これら情報は全て、キーワード付与部 107 によって格納される。

## 【0058】

図 6 に参照回数テーブル 123 のデータ構造の一例を示す。図 6 に示すように、参照回数テーブル 123 は、各収集済み文書が各参照表現で参照されている回数を格納する。これら情報は全て、キーワード付与部 107 によって格納される。

## 【0059】

以下、第 1 実施形態に係わる文書収集装置が実現する特定のコミュニティにとって有用度の高い文書を収集する方法について説明する。説明において以下の表記法を用いる。

・LT (S) は、文書群 S の参照先となる文書群を示す。

・  $X - Y$  は、集合  $X$  と集合  $Y$  の差集合を示す。

#### 【 0 0 6 0 】

最初に、図 7 を用いて特定のコミュニティ向けの文書を収集する処理の大まかな流れについて説明する。まず、収集開始時に、収集済み文書群  $S$  の初期文書群（収集の開始点となる文書群）としてコミュニティ内の文書を与える。

#### 【 0 0 6 1 】

参照関係抽出部 1 0 2 による参照関係の抽出結果及びコミュニティ判別部 1 0 3 による、参照先文書がコミュニティ内の文書であるか否かの判別結果に基づいて、候補判定部 1 0 4 は、次収集候補  $N$  を抽出する（ステップ  $S 1$ ）。次収集候補  $N$  を抽出する処理について、詳しくは後述する。

#### 【 0 0 6 2 】

続いて、文書収集部 1 0 1 は、URL テーブル 1 2 0 に格納された URL に基づいて、次収集候補  $N$  を収集し（ステップ  $S 2$ ）、収集された文書の収集済みフラグをオンにする。これにより、文書収集部 1 0 1 は、新たに収集された次収集候補  $N$  を収集済み文書群  $S$  に加える。つまり、式  $S \cup N$  で示される文書群を新たに収集済み文書群  $S$  とする。

#### 【 0 0 6 3 】

文書収集部 1 0 1 は、収集済み文書群  $S$  に含まれる文書数が規定された文書数以上であるか否かを判定する（ステップ  $S 3$ ）。この判定は、URL テーブル 1 2 0 に格納された収集済みフラグが「オン（1）」になっている文書の数を計数することにより行う。収集済み文書群  $S$  に含まれる文書数が規定された文書数以上でない場合（ステップ  $S 3 : No$ ）、次候補判定部 1 0 4 は、再度次収集候補を決定し（ステップ  $S 4$ ）、ステップ  $S 2$  に戻る。2 回目以降の次収集候補の決定において、今回の収集で新たに収集した文書（以下、新規収集文書という）についての参照関係抽出部 1 0 2 による参照関係の抽出結果、及び、コミュニティ判別部 1 0 3 による新規収集文書の参照先文書がコミュニティ内の文書であるか否かの判別結果に基づいて、候補判定部 1 0 4 は、未収集の参照先文書のうちコミュニティ内の文書を次収集候補  $N$  として抽出する。ステップ  $S 4$  の処理は、ステップ  $S 1$  と同様であるため、ステップ  $S 1$  について、後述する際に一緒に説明す



る。

#### 【 0 0 6 4 】

収集済み文書群 S に含まれる文書数が規定された文書数以上である場合（ステップ S 3 : Y e s）、今度は、候補判定部 1 0 4 は、コミュニティ内外の文書から次収集候補を決定する。そのために、まず、参照関係抽出部 1 0 2 は、新規収集文書の参照関係の抽出し、コミュニティ判別部 1 0 3 は、新規収集文書の参照先文書がコミュニティ内の文書であるか否かを判別する。その後、ランキング部 1 0 5 は、収集済み文書及びその参照先文書、つまり S U L T ( S ) に対して重要度を付与し、重要度に基づいて、未収集の参照先文書、つまり L T ( S ) - S のランキングを行う（ステップ S 5）。このステップ S 5 の処理について詳しくは後述する。

#### 【 0 0 6 5 】

続いて、次候補判定部 1 0 4 は、L T ( S ) - S のうち、コミュニティ内の文書群のランキングで上位 n 1 件に入っている文書群及びコミュニティ外の文書群のランキングで上位 n 2 件に入っている文書群を次収集候補 N とする（ステップ S 6）。このようにコミュニティ内とコミュニティ外とを区別して次収集候補 N を抽出することにより、コミュニティ内またはコミュニティ外に、収集される文書が偏ることを防ぐことができる。

#### 【 0 0 6 6 】

文書収集部 1 0 1 は、URL テーブル 1 2 0 に格納された URL に基づいて、次収集候補 N を収集し（ステップ S 7）、収集された文書の収集済みフラグを「オン（1）」にする。文書収集部 1 0 1 は、URL テーブル 1 2 0 に格納された収集済みフラグが「オン（1）」になっている文書の数を計数することにより、収集済み文書群 S に含まれる文書数が規定された文書数以上であるか否か判定する（ステップ S 8）。

#### 【 0 0 6 7 】

収集済み文書群 S に含まれる文書数が規定された文書数以上でない場合（ステップ S 8 : N o）、ステップ S 5 に戻る。収集済み文書群 S に含まれる文書数が規定された文書数以上である場合（ステップ S 8 : Y e s）、ランキング部 1 0

5、まとめあげ部 1 0 6 及びキーワード部 1 0 7 によって、収集済み文書群 S の文書を選別する（ステップ S 9）。ステップ S 9 の処理について詳しくは後述する。

【 0 0 6 8 】

以下、コミュニティ内の文書を収集する際に、次収集候補を決定する処理について詳しく説明する。この処理は、図 7 のステップ S 1 及びステップ S 4 に相当する。

【 0 0 6 9 】

まず、参照関係抽出部 1 0 2 は、新規収集文書から参照されている参照先文書を抽出する（ステップ S 1 1）。参照関係抽出部 1 0 2 は、各抽出された参照先文書について、参照先文書と同一の URL が URL テーブル 1 2 0 に格納されていない場合、参照先文書の URL を URL テーブル 1 2 0 に格納する（ステップ S 1 2）。同じ URL を重複して格納する必要はないからである。情報を格納する際、参照関係抽出部 1 0 2 は、収集済みフラグを「オフ（0）」とする。

【 0 0 7 0 】

続いて、コミュニティ判別部 1 0 3 は、URL テーブル 1 2 0 に格納された参照先文書の URL の文字列に基づいて、抽出された参照先文書がコミュニティ内の文書であるか否か判別し、コミュニティ内の文書であると判別した場合、コミュニティ判別部 1 0 3 は、URL テーブル 1 2 0 のコミュニティフラグを「オン（1）」とする。それ以外の場合、コミュニティ判別部 1 0 3 は、コミュニティフラグを「オフ（0）」とする（ステップ S 1 3）。さらに、参照関係抽出部 1 0 2 は、コミュニティ判別部 1 0 3 の判別結果に基づいて、参照関係テーブル 1 2 1 の各欄に参照関係を格納する（ステップ S 1 4）。

【 0 0 7 1 】

ここで、本実施形態によれば、コミュニティは、ネットワーク上の文書の集合、つまり文書群として与えられている。従って、同一コミュニティ内の文書であるか否かの判別は、その文書群を示す URL に基づいて判別できる。より具体的には、コミュニティ内の文書であるか否かの判定は、URL の文字列上の特徴に基づいて、以下のようにして行う。

・コミュニティが社内サイトである場合、通常、社内サイトのドメイン名（fujitsu.co.jp等）とドメイン名が同じである文書をコミュニティ内の文書であると判定する。

・コミュニティが業界サイトである場合、その業界サイトに属する複数の企業のサイトのドメイン名のいずれかとドメイン名が同じである文書をコミュニティ内の文書であると判定する。

・コミュニティがユーザグループである場合、各ユーザのサイト（ホームページともいう）のURL（例えば、<http://www.fujitsu.co.jp/foo/>）のいずれかと同じ文字列をURLに含む文書をコミュニティ内の文書であると判定する。

#### 【0072】

次候補判定部104は、収集済み文書の参照先文書であり、かつ、未収集文書である文書LT(S) - Sのうち、コミュニティ内の文書を次収集候補Nとして判定する。具体的には、次候補判定部104は、URLテーブル120を参照し、収集済みフラグが「オフ(0)」であり、且つ、コミュニティフラグが「オン(1)」である文書を次収集候補Nとして決定する（ステップS14）。このような次収集候補Nは、以下の(1)式で表すことができる。

#### 【0073】

$$N = \{d \mid d \in LT(S) - S, d \text{ はコミュニティ内} \} \quad \dots (1)$$

このようにして次収集候補Nを決定し、コミュニティ内の文書をまんべんなく収集することにより、コミュニティ内で必要とされる、意味的に多様な文書についての情報を偏りなく取得することが可能となる。

#### 【0074】

続いて、図9を用いて収集済み文書及びその参照先文書をランキングする処理について説明する。この処理は、図7のステップS5に相当する。

参照関係抽出部102及びコミュニティ判別部103は、新規収集文書の参照関係の抽出し、参照関係をコミュニティの判別結果とともに、URLテーブル120及び参照関係テーブル121に格納する（ステップS21からS23）。このステップS21からS23の処理は、図8で説明したステップS11からS13と同様であるため、詳しい説明は省略する。

## 【 0 0 7 5 】

続いて、ランキング部 1 0 5 は、収集済み文書及びその参照先文書、つまり S U L T ( S ) に対して、参照関係テーブル 1 2 1 に格納された参照関係及び U R L テーブル 1 2 0 に格納された U R L の文字列上の特徴に基づいて重要度を算出し、算出した重要度を U R L テーブル 1 2 0 に格納する（ステップ S 2 4）。ランキング部 1 0 5 は、U R L テーブル 1 2 0 に格納されたコミュニティフラグ及び重要度に基づいて、未収集の参照先文書、つまり、L T ( S ) - S を、コミュニティ内外に分けてランキングする（ステップ S 2 5）。

## 【 0 0 7 6 】

以下、ステップ S 2 4 の重要度を算出する処理について詳しく説明する。上述のように、ランキング部 1 0 5 は、文書の参照関係及び U R L を利用して、収集済み文書の意味内容を分析することなく、文書の重要度を算出する。以下、参照関係に基づいて文書に付与される重要度をリンク重要度という。リンク重要度を付与する際の基本的な考え方は以下の通りである。

- ・類似度の低い U R L から多く参照されている文書は重要である。

## 【 0 0 7 7 】

例えば、一般に、同一サイト内に設けられた複数の文書はそのサイト内の他の文書に参照されているが、それらの文書の U R L は相互に類似する。従って、類似度の高い U R L から参照されている文書の重要度は低いと推定できる。

- ・多くの文書から参照されている文書ほど重要な文書であり、重要な文書から参照されている、U R L の類似度の低い文書は重要である。

## 【 0 0 7 8 】

例えば、有名なディレクトリサービス等及び官公庁等は多くの文書から参照されているが、このような重要な文書から参照されている文書は重要度が高いと考えられる。また、多くの文書やミラーサイトを抱えるサービス（サイト）に設けられた文書等はそのサイト内で参照されていることが多いが、同じサイト内の文書の U R L は大抵類似しているため、「U R L の類似度の低い文書は重要である」という考え方を導入すれば、同じサイトの文書が多く検索されてしまうことを避けることが可能となる。

・URLの類似度は、サーバアドレス、パス、ファイル名の全てが異なるものが最も小さく、ミラーサイトや同一サーバ内の文書は類似度が高くなるように、URLの字面情報から定義する。

【0079】

上述の3つの考え方を導入することにより、全ての参照関係を同等に扱わないでリンク重要度に応じた重みを参照関係に与えることとしている。より具体的には、重みを参照元と参照先文書のURLの類似度の逆数として与えることとしている。以下、リンク重要度の算出についてより詳しく説明する。

【0080】

リンク重要度の算出対象となる文書集合を  $DOC = \{p_1, p_2, \dots, p_n\}$ 、

文書  $p$  のリンク重要度を  $W_p$ 、

文書  $p$  の参照先の文書集合を  $Ref(p)$ 、

文書  $p$  の参照元の文書集合を  $Refed(p)$ 、

文書  $p$  と  $q$  のURL類似度を  $sim(p, q)$ 、

相異度を  $diff(p, q) = 1/sim(p, q)$  とすると、

文書  $p$  から  $q$  に参照が張られているとした時、その参照の重み  $lw(p, q)$  を以下の(1)式で定義する。

【0081】

【数1】

$$\begin{aligned}
 lw(p, q) &= diff(p, q) / \sum_{i \in Ref(p)} diff(p, i) \\
 &= \frac{1}{sim(p, q) \sum_{i \in Ref(p)} \frac{1}{sim(p, i)}} \dots\dots\dots(1)
 \end{aligned}$$

【0082】

この(1)式から分かるように、 $lw(p, q)$  は、 $p$  と  $q$  のURLの類似度  $sim(p, q)$  が低いほど、また、 $p$  からの参照数がより少ないほど大きくなる。

各文書のリンク重要度は、各  $p \in \text{DOC}$  に対して、 $C_q$  を定数（重要度の下限であり、文書によって異なる値を与えてもよい。）として、

【0083】

【数2】

$$W_q = C_q + \sum_{p \in \text{Refed}(q)} W_p * lw(p,q) \quad \dots\dots\dots(2)$$

【0084】

という連立一次方程式の解として定義する。ランキング部105は、この連立一次方程式を解くことにより、リンク重要度を各文書に付与する。なお、連立一次方程式の解法については、既存のアルゴリズムが多数存在するため、説明は省略する。（1）式及び（2）式から、上述の考え方が実現されていることを読み取ることができる。

【0085】

次に、（1）式及び（2）式中の文書  $p$  と  $q$  の URL 類似度  $\text{sim}(p,q)$  について説明する。URL 類似度は、ランキング部105のURL判別部（不図示）により算出される。一般に、文書のURLは、サーバアドレス、パス、ファイル名の三種類の情報から構成される。例えば、WWW文書のURL、

<http://www.flab.fujitsu.co.jp/hypertext/news/1999/product1.html> は、サーバアドレス（[www.flab.fujitsu.co.jp](http://www.flab.fujitsu.co.jp)）、パス（[hypertext/news/1999](http://www.flab.fujitsu.co.jp/hypertext/news/1999)）、ファイル名（[product1.html](http://www.flab.fujitsu.co.jp/hypertext/news/1999/product1.html)）の3種類の情報から構成される。

【0086】

本実施形態では、与えられた2つの文書  $p$  及び  $q$  のURL類似度を、上記の三種類の組合せにより定義する。類似度  $\text{sim}(p,q)$  として、例えば、以下に述べるドメイン類似度  $\text{sim\_domain}(p,q)$  及び融合類似度  $\text{sim\_merge}(p,q)$  が考えられる。

【0087】

ドメイン類似度  $\text{sim\_domain}(p,q)$  は、ドメインの類似に基づいて算出される

。ドメインとは、サーバアドレスの後半部分であり、会社や組織を表す。サーバアドレスが.com、.edu、.org等で終わる米国サーバの場合はサーバアドレスの後ろから2つめまで、サーバアドレスが.jp、.fr等で終わる他国のサーバの場合はサーバアドレスの後ろから3つめまでがドメインに相当する。

## 【 0 0 8 8 】

文書  $p$  と文書  $q$  のドメイン類似度は以下の式により定義される。

$$\begin{aligned} \text{sim\_domain}(p, q) &= 1 / \alpha && (p, q \text{ が同一ドメインの場合}) \\ &= 1 && (p, q \text{ が異なるドメインの場合}) \end{aligned}$$

ここで、 $\alpha$  は定数で、0 より大きく 1 より小さい実数値を取るとする。

## 【 0 0 8 9 】

また、 $\text{sim}(p, q)$  として、前述の三種類の情報を融合した類似度  $\text{sim\_merge}(p, q)$  を次のように定義する。

$\text{sim\_merge}(p, q) = (\text{サーバアドレスの類似度}) + (\text{パスの類似度}) + (\text{ファイル名の類似度})$

以下、右辺の各項の算出方法について説明する。

## 【 0 0 9 0 】

サーバアドレスの類似度は、アドレスの階層を後ろから見ていき、 $n$  レベルまで一致した場合、類似度を  $1 + n$  とする。例えば、`www.fujitsu.co.jp` と `www.flab.fujitsu.co.jp` は 3 レベルまで一致しているので 4 となる。`www.fujitsu.co.jp` と `www.fujitsu.com` は 1 レベルも一致していないので（一致 0 レベル）、類似度は 1 である。

## 【 0 0 9 1 】

パスの類似度は、先頭からパスの"/" で区切られた要素毎に比較し、一致したレベルまでを類似度とする。例えば、`/doc/patent/index.html` と `/doc/patent/1999/2/file.html` とは、2 レベルまで一致しているので類似度は 2 である。

## 【 0 0 9 2 】

ファイル名の類似度は、ファイル名が一致する場合、類似度 1 とする。

この  $\text{sim\_merge}(p, q)$  によっても、URL が似通った文書が多く検索されることを防ぐことができる。

## 【0093】

このようにして、ランキング部105は、文書に重要度を付与し、高い重要度を付与された文書を上位にランキングする。

このように、本実施形態によれば、ランキング部105は、取得した文書の参照関係及びURLの文字列の特徴に基づいて、文書本文の意味内容を解析せずに、つまり処理速度が速くかつ精度良く、文書に重要度を付与し、その重要度に基づいて文書をランキングすることができる。

## 【0094】

以下、図10を用いて収集済み文書を選別する処理について詳しく説明する。この処理は図7のステップS9に相当する。まず、まとめあげ部106は、収集済み文書群Sで用いられている参照表現に基づいて、収集済み文書群Sをまとめあげる（ステップS31）。なお、参照表現とは、例えば、HTML（Hyper Text Mark-up Language）では、アンカータグで囲まれた部分がそれに相当する。

## 【0095】

より具体的には、予め不図示のまとめあげ参照表現テーブルに、「次に」、「前へ」といった参照表現（参照時に用いられる文字列）を格納する。これら「次に」、「前へ」といった参照表現を用いている文書は、同一内容であるが、URLが分散されている文書と推定される。まとめあげ部106は、まとめあげ参照表現テーブルに格納されている参照表現を文書から抽出し、以下のようにして文書をまとめあげる。

- ・文書doc1の中から「次へ」、「次に続く」、「Next」というような表現により、文書doc2が参照されている場合、まとめあげ部106は、文書doc2を文書doc1に縮退する。この操作の繰り返しを可能な限り行う。

- ・文書doc1の中から「前へ」、「前に戻る」、「Prev」といった表現により、文書doc2が参照されている場合、まとめあげ部106は、文書doc1をdoc2に縮退する。この操作の繰り返しを可能な限り行う。

## 【0096】

続いて、キーワード付与部107は、参照表現に基づいて収集済み文書Sにキーワードを付す（ステップS32）。キーワード付与処理について詳しくは後述



する。最後に、ランキング部105は、上述の図9のステップS24と同様にし  
て、収集済み文書に重要度を付与し、重要度をURLテーブル120に格納する  
。ランキング部105は、重要度に基づいて収集済み文書をランキングする（ス  
テップS33）。

#### 【0097】

次に、ステップS32のキーワード付与処理について、図11を用いて詳しく  
説明する。まず、予め、収集済み文書で用いられている参照表現のうち、「ホー  
ムへ」、「トップに戻る」等、参照先文書に関係なく、しばしば使用される参照  
表現を不図示の不要語辞書に格納する（不図示）。キーワード付与部107は、  
収集済み文書群Sから参照表現を抽出し、各参照表現wについて、参照表現wを  
用いて参照される相異なる文書の数DF(w)を集計し、参照表現wを識別する  
表現ID、その参照表現（文字列）とともにDF(w)の集計結果を参照表現テ  
ーブル122に格納する（ステップS41）。この段階では、要否フラグを「オ  
フ(0)」としておく。

#### 【0098】

キーワード付与部107は、参照表現wのうち、DF(w)が所定の数以上で  
あるものをキーワード候補から省く（ステップS42）。言い換えると、参照先  
文書まで含めた総文書数をNとすると、以下の(3)式に該当する参照表現wを  
省く。

#### 【0099】

$$DF(w) > \alpha N \quad \dots\dots (3)$$

ここで、 $\alpha$ は、定数であり、例えば0.1としてもよい。

キーワード付与部107は、参照表現wのうち、不要語辞書に格納されている  
特定の参照表現をキーワード候補から省く（ステップS43）。これらの参照表  
現は、参照先文書に関係なく使用されているため、キーワードとして用いるには  
適切でないからである。

#### 【0100】

キーワード付与部107は、収集済み文書Sから、文書dを取り出し、収集済  
み文書群Sとdの差集合、つまりS-dを新たな収集済み文書群Sとする（ステ

ップ S44)。

#### 【0101】

キーワード付与部 107 は、キーワード付与部 107 は、文書 d で各参照表現 w が参照されている回数  $TF(d, w)$  を集計し、以下の (4) 式を用いて、文書 d について各参照表現 w の重み  $W(d, w)$  を算出する (ステップ S45)。

#### 【0102】

$$W(d, w) = TF(d, w) \log(N/DF(w)) \dots (4)$$

キーワード付与部 107 は、参照表現テーブル 122 にアクセスし、参照表現の重み w の大きい順に高々 n 個の参照表現の要否フラグを「オン (1)」とする。つまり、重み w の大きい順に高々 n 個の参照表現を文書 d のキーワードとする。

#### 【0103】

このようにして得られた参照表現に基づくキーワードは、文書 d の本文に含まれる単語に基づくキーワードと異なり、様々な異称をキーワードとして取得することが特徴の 1 つである。例えば、ある企業のホームページへの参照表現から、その企業の様々な呼称 (正式名、略称、通称、英語名) を取得することができる。また、例えば、用語「Linux」に関して、「リナックス」、「ライナックス」等の様々な内証がキーワードとして取得することができる。一方、一般に 1 つの文書の本文ではこうした異称のうち 1 つだけを統一的に用いるため、本文からキーワードを取得する場合では異称をキーワードとして取得することはできない。

#### 【0104】

また、参照表現から取得したキーワードに、文書 d の本文に出現する単語のうちで頻出する単語からキーワード及び文書 d を示す URL から得たキーワード、例えば、<http://www.fujitsu.com/> であれば、キーワードとして fujitsu、を加えることとしてもよい。これにより、文書 d に多様なキーワードを付与することが可能になる。

#### 【0105】

図 12 に、第 1 実施形態に係わる文書収集装置を用いて収集した文書をユーザ

に提供する画面の一例を示す。図12において、収集した優良コンテンツ130を、分類エンジン150を用いてディレクトリに分け、サーバ160のクライアントに提供する場合を例としている。クライアントは、画面180でキーワードを入力する、又は、カテゴリを選択することにより、閲覧したい文書へのリンクまたはリンク集を画面に表示させることができる。

#### 【0106】

クライアントがキーワードを入力した場合、画面181に示すようにキーワードに基づいて検索された文書へのリンクが、重要度と共に表示される。本実施形態によれば、入力されたキーワードの異称も合わせて検索することが可能である。カテゴリを選択した場合、画面182に示すように選択されたカテゴリに関連する文書へのリンク集が表示される。

#### 【0107】

ここで、画面181及び画面182に示すように、検索された文書を提示する際に、URLテーブル120に格納されたコミュニティフラグに基づいて、文書をコミュニティ内外に分けて提示することとしても良い。

#### 【0108】

以下、第2実施形態に係わる文書収集装置について説明する。第2実施形態に係わる文書収集装置は、特定分野に関する文書を収集する。本実施形態に係わる文書収集装置において以下の考え方を採用する。

・ネットワークにおいて、参照の親子／兄弟関係にある文書は、内容的に似通っている傾向にある。ある程度の文書群としばしば親子／兄弟関係にあるとされる文書は、元文書群と同じような分野の内容である可能性が高い。元の文書群からと親子／兄弟関係にある文書のうち参照度（親子関係）や共参照度（兄弟関係）の高い文書を収集し、元文書群に繰り込み、という操作を多段階に繰り返すことで、当該分野に関する文書を収集していくことができる。

#### 【0109】

図13に第2実施形態に係わる文書収集装置の構成を示す。図13に示すように第2実施形態に係わる文書収集装置200は、文書収集部101、参照関係抽出部102、候補判定部104、参照度／共参照度算出部201、ランキング部

105、まとめあげ部106及びキーワード付与部107を備える。参照度／共参照度算出部201は、文書の参照関係に基づいて、ある文書が特定分野に関連している度合いを算出する。その他の各部の機能は、第1実施形態で説明した通りである。

#### 【0110】

第2実施形態に係わる文書収集装置において、収集開始に先立って、まず、当該分野の代表的な文書を既存の検索エンジンやリンク集を用いて収集し、正例文書群PSとして与え、当該分野と重ならない任意の文やの文書も同様にして収集して負例文書群NSとして与え、PSUNSを収集済み文書群Sとする。この収集済み文書群Sが収集の開始点となる。

#### 【0111】

参照関係抽出部102は、収集済み文書群Sから参照関係を抽出し、収集済み文書群Sの参照先文書のURLをURLテーブル120に格納し、抽出された参照関係を参照関係テーブル121に格納する。ここで、第2実施形態に係わる文書収集装置において、URLテーブル120に、コミュニティフラグの代わりに正例文書群PSに含まれる文書であるか否かを示す正例フラグの欄を含む。正例フラグは、正例文書群PSに含まれる文書である場合に「オン(1)」となる。また、参照関係テーブル121に参照関係を格納する際、コミュニティ内外で分けることは不要となる。

#### 【0112】

参照度／共参照度算出部201は、参照関係抽出部102が抽出した参照関係に基づいて、正例文書群PS及び負例文書群NSと収集済み文書Sの参照先文書との関係を示す参照度及び共参照度を算出する。次候補判定部104は、参照度／共参照度算出部201が算出した参照度及び共参照度に基づいて、収集済み文書群Sの参照先文書であって、正例文書群PSに含まれない文書のなかから所定の条件を満たす文書を次収集候補Nとして判定する。次候補判定部104は、次収集候補Nのうち負例文書群NSに含まれている文書を負例文書群NSから除き、正例文書群PSに加える。

#### 【0113】

文書収集部 1 0 1 は、URL テーブル 1 2 0 を参照し、次収集候補 N のうち未収集文書を収集し、収集した文書を正例文書群 P S に加える。第 2 実施形態に係わる文書収集装置 2 0 0 は、正例文書群 P S の文書数が規定された数以上になるまで、上述のようにして収集済み文書 S の参照関係を抽出し、参照関係に基づいて次収集候補 N を決定し、次収集候補 N を収集する処理を繰り返す。

## 【 0 1 1 4 】

収集済み文書 S が規定された数以上になると、まとめあげ部 1 0 6 は参照表現に基づいて収集済み文書群 S をまとめあげ、キーワード付与部 1 0 7 は参照表現が用いられる頻度等に基づいて収集済み文書群 S にキーワードを付す。ランキング部 1 0 5 は、参照関係及び URL の文字列上の特徴に基づいて各収集済み文書 S の重要度を算出し、重要度に基づいて収集済み文書 S をランキングする。これにより、分野別優良コンテンツ 2 1 0 を作成する。このように、第 2 実施形態に係わる文書収集装置によれば、文書本文の内容を解析せずに、特定分野に関する文書を収集し、まとめあげ、キーワードを付与することができる。

## 【 0 1 1 5 】

分野別優良コンテンツ 2 1 0 は、検索エンジン 1 4 0 を介してサーバ 1 6 0 に提供される。サーバのクライアントはブラウザ 1 6 0 を用いて検索サービスの提供を受けることができる。

## 【 0 1 1 6 】

以下、第 2 実施形態に係わる文書収集装置が実現する特定分野に関する文書収集方法について説明する。まず、用いる表記法について説明する。

- ・ L T ( B ) は、文書群 B の参照先文書集合を示す。
- ・ L T ( p ) は、文書 p の参照先文書集合を示す。
- ・ L S ( d , X ) = { c ∈ X | c refers d } は、文書集合 X のうち文書 d を参照している文書の集合を示す。
- ・ L S ( A , X ) = { c ∈ X | ∃ d ∈ A , c refers d } は、文書集合 X のうち集合 A 中の少なくとも 1 文書を参照している文書の集合を示す。
- ・ C C ( d , A , X ) = L S ( d , X ) ∩ L S ( A , X ) は、文書集合 X のうちで、文書 d、及び集合 A の文書（少なくとも 1 文書）の両方を参照している文書

の集合を示す。

【0117】

図14に、LT(S)、LT(p)、LS(d, X)及びLS(A, X)について、各集合が意味する文書の参照関係を示す。図14において黒丸は文書を示し、矢印は参照関係を示し、矢印の元が参照元、矢印の先が参照先を示す。図14に示すように、LT(B)とLS(A, X)及びLT(p)とLS(d, X)は、それぞれ矢印が逆になっている、つまり参照先文書と参照元文書が入れかわった関係にあることが分かる。また、図15に、CC(d, A, X)が意味する文書の参照関係を示す。

【0118】

以下、図16を用いて特定分野に関する文書を収集する処理について説明する。第2実施形態に係わる文書収集装置によれば、「XML」や「Linux」といった、特定分野（ジャンル）に関する意味的に類似した文書を優先的に収集する場合に、文書本文の内容を解析する処理を行わずに、参照関係に基づいて収集することが可能である。

【0119】

まず、当該分野に属する代表的な文書を、既存の検索エンジンやリンク集から探し出して収集し、正例文書群PSとする。同様にして当該分野とは重ならない分野に属する文書を、探し出して収集し、負例文書群NSとする。この正例文書群PSと負例文書群NSが初期文書群となる。そして、PS及びNSの文書のURL、収集済みフラグ（全て「オン（1）」）、及び正例フラグ（正例文書の場合「オン（1）」）をURLテーブル120に格納する。正例文書群PSと負例文書群NSの和集合PS∪NSを収集済み文書群Sとする（ステップS51）。ここで、例えば、当該分野を「コンピュータ」であるとする、当該分野と重ならない分野の例として、「手芸」、「料理」、「美容」等が考えられる。

【0120】

参照関係抽出部102は、収集開始時は初期の収集済み文書群S（初期文書群）から、それ以降は新規収集文書から参照関係を抽出し（ステップS52）、参照先文書のURLをURLテーブル120に格納し、参照関係を参照関係テーブ

ル 1 2 1 に格納する。この処理は、第 1 実施形態と同様である。

### 【0121】

参照度／共参照度算出部 201 は、抽出された参照関係に基づいて、収集済み文書群 S の参照先文書から正例文書群 P S に含まれる文書を除いた文書集合 T ( S ) = L T ( S ) - P S に含まれる文書  $d \in T ( S )$  について、以下の (5) 式を用いて参照度  $R_{\text{score}} (d, P S, S)$  を算出する。次候補判定部 105 は、参照度  $R_{\text{score}} (d, P S, S)$  が上位 n 1 件に入っている文書群を N 1 とする。(ステップ S 53)。なお、収集済み文書が正例文書群 P S に含まれるか否かは、URL テーブル 120 の正例フラグを参照することにより判定できる。

### 【0122】

【数 3】

$$R_{\text{score}}(d, P S, S) = \log(|L S(d, P S)|) \cdot \frac{|L S(d, P S)|}{|L S(d, S)|}$$

### 【0123】

(5) 式の第 1 項は、文書 d を参照している正例文書群 P S の文書数の対数を示す。また、(5) 式の第 2 項は、文書 d を参照している収集済み文書数に対する、文書 d を参照している正例文書群 P S の文書数の割合を示す。従って、S のうち正例文書群 P S からのみ多く参照されている文書 d ほど、 $R_{\text{score}} (d, P S, S)$  が大きな値を取ることが分かる。

### 【0124】

つまり、次候補判定部 105 は、参照度  $R_{\text{score}} (d, P S, S)$  に基づいて、新規収集文書の参照先文書のうち、特定分野に関係ある正例文書群 P S から多く参照され、特定分野とあまり関係ない負例文書群 N S から参照されていない文書を N 1 として決定する。図 17 に、文書 d について参照度を算出する際に、(6) 式に含まれる各集合が意味する参照関係を示す。

### 【0125】

続いて、参照度／共参照度算出部 201 は、文書  $d \in T ( S ) - N 1$  について

、以下の(6)式を用いて共参照度  $C_{\text{score}}(d, PS, S)$  を算出する。次候補判定部105は、 $d \in T(S) - N1$  のうちで共参照度  $C_{\text{score}}(d, PS, S)$  が上位  $n2$  件に入っている文書群を  $N2$  とする(ステップS54)。

【0126】

【数4】

$$C_{\text{score}}(d, PS, S) = \log \left( \sum_{p \in CC(d, s)} |LT(p) \cap PS| \right) \frac{|CC(d, PS, S)|}{|LS(d, S)|}$$

【0127】

(6)式の第1項の対数の中身は、文書  $d$  及び正例文書群  $PS$  の文書の両方を参照している収集済み文書  $p$  全てについての、文書  $p$  の参照先文書であって正例文書群  $PS$  に含まれる文書数の積和を示す。従って、共参照度  $C_{\text{score}}(d, PS, S)$  は、文書  $d$  及び正例文書群  $PS$  の文書の両方を参照している収集済み文書  $p$  の数が多い文書  $d$  ほど、及び、このような文書  $p$  の参照先文書であって正例文書群  $PS$  に含まれる文書の数が多いような文書  $d$  ほど、大きな値を取ることが分かる。言い換えると、正例文書群  $PS$  の文書を参照している収集済み文書から参照されている文書  $d$  について、その文書  $d$  を参照している収集済み文書の数が多い文書  $d$  ほど、共参照度  $C_{\text{score}}(d, PS, S)$  は、大きな値を取る。

【0128】

(6)式の第2項は、文書  $d$  の参照元となっている収集済み文書の数に対する、文書  $d$  と共に参照されている文書  $p$  の数の割合を示す。共参照度  $C_{\text{score}}(d, PS, S)$  は、この割合が大きいほど大きな値を取る。図18に、文書  $d$  について共参照度を算出する際に、(6)式に含まれる各集合が意味する参照関係を示す。

【0129】

次候補判定部105は次収集候補  $N = N1 \cup N2$  とする(ステップS55)。次候補判定部105は、次収集候補  $N$  のURLをキーとしてURLテーブル120を検索し、次収集候補  $N$  の正例フラグを「オン(1)」する。この処理により



、負例文書群 N S に含まれていたが、次収集候補として判定された文書が、負例文書群 N S から除かれ、正例文書群 P S に加えられることとなる（ステップ S 5 6）。

#### 【0 1 3 0】

文書収集部 1 0 1 は、URL テーブル 1 2 0 に格納された URL に基づいて、次収集候補 N のうち未収集文書をネットワークから収集し、収集した文書に対応する収集済みフラグを「オン（1）」にする（ステップ S 5 7）。この処理により、新規収集文書を正例文書群 P S に加える。文書収集部 1 0 1 は、URL テーブル 1 2 0 を参照し、正例文書群 P S の文書数が規定された数以上であるか否か判定する（ステップ S 5 8）。正例文書群 P S の文書数が規定された数以上でない場合（ステップ S 5 8 : N o）、ステップ S 5 2 に戻って処理を繰り返す。

#### 【0 1 3 1】

正例文書群 P S の文書数が規定された数以上である場合（ステップ S 5 8 : Y e s）、正例文書群 P S の文書を選別し（ステップ S 5 9）、処理を終了する。文書を選別処理は、第 1 実施形態と同様であるため説明を省略する。

#### 【0 1 3 2】

このようにして、本実施形態によれば、文書本文の内容を解析することなく、特定分野に関する文書を精度よく、かつ迅速に収集することが可能となる。

以下、第 2 実施形態の変形例について説明する。負例文書群 N S は、集めることも難しいため、収集処理の後に廃棄することをさせて、有効利用することが望ましい。そこで、第 2 実施形態の変形例に係わる文書収集装置によれば、上記処理で収集した負例文書群 N S を有効に利用することとする。これにより、なるべく独立な、例えば、「J a v a 言語」と「編物」及び「フランス料理」等、複数分野の文書を並行して収集することを可能とする。そのために、ある分野の文書を収集する際、その分野の文書群を正例文書群とし、その分野以外の他の分野の文書群を負例文書群 N S として扱う。

#### 【0 1 3 3】

文書収集装置の構成は、図 1 3 を用いて説明した通りであるため、説明を省略する。以下、図 1 9 を用いて第 2 実施形態の変形例に係わる文書収集装置で行う

処理について説明する。

【0134】

まず、 $n$  個の独立な分野の文書群  $D_i$  ( $i = 1, 2, \dots, n$ ) を、検索エンジンやリンク集等から探し出して収集し、文書群  $D_i$  の文書の URL、収集済みフラグ、及び分野を識別する情報である分野識別情報を URL テーブル 120 に格納する。第 2 実施形態に係わる文書収集装置では、正例フラグは不要である。文書群  $D_i$  は、分野  $i$  の初期文書群となる。収集済み文書群を  $D = (D_1, D_2, \dots, D_n)$  とする (ステップ S61)。

【0135】

まず、参照関係抽出部 201 は、 $i$  を与える (ステップ S62)。なお、収集開始時に、参照関係抽出部 201 は、 $i$  を 1 とする。続いて、参照関係抽出部 201 は、 $i$  が  $n$  を超えているか否か判定する (ステップ S63)。 $i$  が  $n$  を超えている場合 (ステップ S63: Yes)、ステップ S71 に進む。そうでない場合 (ステップ S63: No)、参照関係抽出部 102 は、分野  $i$  に対応する文書群  $D_i$  の新規収集文書から (収集開始時は初期文書群から)、参照関係を抽出し、参照先文書の URL を URL テーブル 120 に、参照関係を参照関係テーブル 121 にそれぞれ格納する (ステップ S64)。この処理は、第 1 実施形態と同様である。

【0136】

参照度／共参照度算出部 201 は、文書群  $D_i$  の参照先文書であって、収集済み文書群  $D$  に含まれない文書群  $T(D_i) = LT(D_i) - D$  を次収集範囲とし、この次収集範囲  $T(D_i)$  に含まれる文書  $d \in T(D_i)$  について、上述の (5) 式を用いて参照度  $R_{\text{score}}(d, D_i, D)$  を算出する。次候補判定部 105 は、参照度  $R_{\text{score}}(d, D_i, D)$  が上位  $n1$  件に入っている文書群を  $N1_i$  とする。 (ステップ S65)。なお、収集済み文書が含まれる分野は、URL テーブル 120 の分野識別情報を参照することにより判定できる。

【0137】

参照度／共参照度算出部 201 は、次収集範囲  $T(D_i)$  から  $N1_i$  を除いた集合に含まれる文書  $d \in T(D_i) - N1_i$  について、上述の (6) 式を用いて共

参照度  $C_{\text{score}}(d, D_i, D)$  を算出する。次候補判定部 105 は、共参照度  $C_{\text{score}}(d, D_i, D)$  が上位  $n$  件に入っている文書群を  $N_{2i}$  とする。(ステップ S66)。

【0138】

次候補判定部 105 は、 $N_{1i} \cup N_{2i}$  を分野  $i$  についての次収集候補  $N_i$  とする(ステップ S67)。次候補判定部 105 は、URL テーブル 120 にアクセスし、次収集候補  $N_i$  に現在の  $i$  の値に対応した分類識別情報を付す。文書収集部 101 は、ネットワークから次収集候補  $N_i$  を収集する(ステップ S68)。文書収集部 101 は、URL テーブル 120 にアクセスし、収集された次収集候補  $N_i$  (新規収集文書群)の収集済みフラグを「オン(1)」とする。これにより、文書収集部 101 は、文書群  $D_i$  に新規収集文書群を加えて新たな文書群  $D_i$  とする(ステップ S69)。

【0139】

続いて、参照関係抽出部 102 は、 $i$  を 1 インクリメントし(ステップ S70)、ステップ S63 に戻る。文書収集装置 200 は、上述の処理を  $i$  が  $n$  を超えるまで、処理を繰り返す。

【0140】

$i$  が  $n$  を超えると(ステップ S63: Yes)、参照関係抽出部 102 は、URL テーブル 120 を参照し、収集済みフラグ及び分野識別情報に基づいて、各文書群  $D_i$  の文書数を計数し、各文書群  $D_i$  の文書数が規定された数以上であるか否か判定する(ステップ S71)。文書数が規定数以上でない文書群  $D_k$  ( $k$  は 1 から  $n$  までの任意の数)がある場合、ステップ S62 に戻り、参照関係抽出部 102 は、 $i = k$  としてステップ S63 以下の処理を繰り返す。

【0141】

なお、文書数が規定数以上でない文書群  $D_k$  が複数ある場合、例えば、 $D_{k1}$ 、 $D_{k2}$  及び  $D_{k3}$  がある場合、 $i = k1$ 、 $k2$  及び  $k3$  である場合について、ステップ S63 以下の処理を繰り返す。 $D_1$  から  $D_n$  まで全ての収集済み文書群  $D_i$  について文書数が規定数以上である場合(ステップ S71: Yes)、処理を終了する。

## 【0142】

これにより、ある分野の文書を収集する際に、その分野の文書群を正例文書群  $PS$  とし、他の残りの分野の文書群の和集合を負例文書群  $NS$  として用いることができるため、負例文書群  $NS$  に関する処理が無駄にならないこととなる。

## 【0143】

また、第2実施形態の変形例によれば、ある分野の文書群  $D_i$  を正例文書群  $PS$  として、その分野に関する文書を収集する場合に注目すると、負例文書群  $NS$  として用いられる他の分野の文書群が、正例文書群  $PS$  と比べ大きくなる。さらにまた、負例文書群  $NS$  自体も他の分野に関する文書群であるため、意味的に一定している。変形例ではない第2実施形態においてある程度以上収集が進むと、正例文書群  $PS$  が大きくなる一方で負例文書群  $NS$  から正例文書群  $PS$  に文書が移されることによって、例えば(5)式に示される  $R_{score}(d, PS, S)$  の第2項が大きくなっていく事態が生じうる。これによって、収集の精度が低下する可能性があったが、変形例ではその可能性が低くなる。

## 【0144】

以下、図20及び図21を用いて、第2実施形態に係わる文書収集装置において特定分野に関する文書を収集する精度について説明する。図20に、ネットワークから収集した約670万URLの文書を全体集合  $D$  とし、URLに「Linux」を含む15,000URLの文書を正解例  $L$  とし、任意に選択した約5,000URLの文書を正例文書群  $PS$  の初期文書とし、及びそれ以外の文書 ( $D - PS$ ) を負例文書群  $NS$  の初期文書として、文書収集装置の収集精度を実験した結果を示す。

## 【0145】

図20において、横軸に収集のくり返し回数  $i$ 、縦軸に適合率又は再現率を示す。再現率を折れ線、適合率を四角プロットで示す。ここで、 $i$  回目の繰り返しで得られた正例集合  $PS_i$  についての適合率及び再現率は、以下(7)式及び(8)式で示される。

## 【0146】

$$\text{適合率} = |PS_i \cap L| / |PS_i| \cdots (7)$$

$$\text{再現率} = |PS_i \cap L| / |L| \cdots (8)$$

つまり、適合率は、正例集合  $S_i$  中の正例文書群  $S$  に含まれる正解例  $L$  の割合であり、対象としている分野に含まれない文書（いわゆるゴミ）の少なさを示す。再現率は、正解例  $L$  中の正例文書群  $S_i$  に含まれる正解例  $L$  の割合であり、対象としている分野に含まれる文書が収集されないこと（いわゆる漏れ）の少なさを示す。図 20 に示すように、繰り返し回数が 73 回程度になると、再現率が急激に低下するが、数十回の繰り返しでは、適合率、再現率とも良好であることが分かる。なお、繰り返し回数が 73 回程度になると再現率が低下する原因は、所謂ゴミがゴミをよぶためであると考えられる。

【0147】

図 21 に、URL に「What's New」を含む 14,000 URL を正解例  $L$  とした場合に、同様の実験を行った結果を示す。図 21 に示すように、繰り返し回数が数回程度になると急激に適合率が低下している。これは、What's New のようなコンテンツは、互いにあまり意味的な関連（つながり）が無いためと考えられる。

【0148】

図 20 に示す実験結果から、本実施形態に係わる文書収集装置によれば意味的に関連する文書群を効率よく収集することができることが分かる。

上述において説明した各サーバ及び各端末は、図 22 に示すような情報処理装置（コンピュータ）を用いて構成することができる。図 22 の情報処理装置 20 は、CPU 301、メモリ 302、入力装置 303、出力装置 304、外部記憶装置 305、媒体駆動装置 306、及びネットワーク接続装置 307 を備え、それらはバス 308 により互いに接続されている。

【0149】

メモリ 302 は、例えば、ROM (Read Only Memory)、RAM (Random Access Memory) 等を含み、処理に用いられるプログラムとデータを格納する。CPU 301 は、メモリ 302 を利用してプログラムを実行することにより、必要な処理を行う。

【0150】

上述の各サーバ及び各端末を構成する各機器及び各部は、それぞれメモリ 30

2の特定のプログラムコードセグメントにプログラムとして格納される。入力装置303は、例えば、キーボード、ポインティングデバイス、タッチパネル等であり、ユーザからの指示や情報の入力に用いられる。出力装置304は、例えば、ディスプレイやプリンタ等であり、情報処理装置300の利用者への問い合わせ、処理結果等の出力に用いられる。

#### 【0151】

外部記憶装置305は、例えば、磁気ディスク装置、光ディスク装置、光磁気ディスク装置等である。この外部記憶装置305に上述のプログラムとデータを保存しておき、必要に応じて、それらをメモリ302にロードして使用することもできる。

#### 【0152】

媒体駆動装置306は、可搬記録媒体309を駆動し、その記録内容にアクセスする。可搬記録媒体309としては、メモリカード、メモリスティック、フロッピーディスク、CD-ROM (Compact Disc Read Only Memory)、光ディスク、光磁気ディスク、DVD (Digital Versatile Disk) 等、任意の情報処理装置で読み取り可能な記録媒体が用いられる。この可搬記録媒体309に上述のプログラムとデータを格納しておき、必要に応じて、それらをメモリ302にロードして使用することもできる。

#### 【0153】

ネットワーク接続装置307は、LAN、WAN等の任意のネットワーク（回線）を介して外部の装置と通信し、通信に伴うデータ変換を行う。また、必要に応じて、上述のプログラムとデータを外部の装置から受け取り、それらをメモリ302にロードして使用することもできる。

#### 【0154】

図23は、図22の情報処理装置300にプログラムとデータを供給することのできる情報処理装置で読み取り可能な記録媒体及び伝送信号を示している。

なお、本発明は、情報処理装置により使用されたときに、上述の本発明の実施形態の各構成によって実現される機能と同様の機能を情報処理装置に行わせるための情報処理装置で読み出し可能な記録媒体309として構成することもできる

## 【 0 1 5 5 】

実施形態において各装置により行なわれる処理と同様のものを情報処理装置に行なわせるプログラムを、情報処理装置で読み取り可能な記録媒体 3 0 9 に予め記憶させておき、図 2 3 に示すようにしてその記録媒体 3 0 9 からそのプログラムを情報処理装置 3 0 0 に読み出させてその情報処理装置 3 0 0 のメモリ 3 0 2 や外部記憶装置 3 0 5 に一旦格納させ、その情報処理装置 3 0 0 の有する CPU 3 0 1 にこの格納されたプログラムを読み出させて実行させる。

## 【 0 1 5 6 】

また、プログラム（データ）提供者 3 1 0 から情報処理装置 3 0 0 にプログラムをダウンロードする際に回線 3 1 1（伝送媒体）を介して伝送される伝送信号自体も、上述した本発明の実施形態において説明した各装置に相当する機能を汎用的な情報処理装置で行なわせることのできるものである。

## 【 0 1 5 7 】

以上、本発明の実施形態について説明したが、本発明は上述した実施形態に限定されるものではなく、他の様々な変更が可能である。

例えば、第 1 実施形態に係わる文書収集装置 1 0 0 と第 2 実施形態に係わる文書収集装置 2 0 0 とを組みせるように構成することにより、コミュニティ向けに分野別に文書を収集させることとしてもよい。

## 【 0 1 5 8 】

また、文書収集装置 1 0 0 又は 2 0 0 を構成する各部及び各 DB は、お互いに連携して動作することにより一連のビジネスプロセスを実現する。これら各部及び各 DB は同じサーバに設けられてもよいし、異なるサーバに設けられネットワークを介して連携して動作することとしてもよい。

## 【 0 1 5 9 】

（付記 1） ネットワークから文書を収集する文書収集方法であって、  
前記文書の参照関係に基づいて、前記ネットワーク上のコミュニティ内から文書を所定数以上収集し、

前記コミュニティから前記所定数以上の文書を収集した後、収集済み文書の参

照関係に基づいて、前記コミュニティ内外から文書を収集する、  
ことを特徴とする文書収集方法。

【 0 1 6 0 】

(付記 2) 前記収集済み文書群の参照関係及びネットワーク上の場所を示す  
情報に基づいて重要度を算出し、  
前記参照関係及び前記重要度に基づいて、収集すべき文書を決定する、  
ことを特徴とする付記 1 記載の文書収集方法。

【 0 1 6 1 】

(付記 3) 前記収集すべき文書は、前記コミュニティ内外別に決定される、  
ことを特徴とする付記 2 記載の文書収集方法。

(付記 4) 前記収集済み文書群を検索した結果を、前記コミュニティ内外に  
分けて提示する、  
ことを特徴とする付記 3 記載の文書収集方法。

【 0 1 6 2 】

(付記 5) 前記コミュニティ内の文書であるか否かを前記ネットワーク上の  
場所を示す情報に基づいて判定する、  
ことを特徴とする付記 2 記載の文書収集方法。

【 0 1 6 3 】

(付記 6) ネットワークから文書を収集する文書収集方法であって、  
ある分野に関する文書群である正例文書群と、前記分野と関連が少ない分野に  
関する文書群である負例文書群とを与え、

前記正例文書群及び前記負例文書群の参照関係に基づいて、前記分野に関する  
収集すべき文書を決定し、 前記ネットワークから前記収集すべき文書を収集す  
る、

ことを特徴とする文書収集方法。

【 0 1 6 4 】

(付記 7) 前記参照関係に基づいて、前記正例文書群の文書からのみ参照さ  
れる度合いである参照度を算出し、

前記参照度が高い文書を収集すべき文書として決定する、



ことを特徴とする付記 6 記載の文書収集方法。

【 0 1 6 5 】

(付記 8) 前記参照関係に基づいて、前記正例文書群の文書を参照している収集済み文書から参照されている文書について、収集済み文書からの被参照数を示す共参照度を算出し、

共参照度が高い文書を収集すべき文書として決定する、

ことを特徴とする付記 6 記載の文書収集方法。

【 0 1 6 6 】

(付記 9) 前記負例文書群は、複数の分野に関する文書群の和集合である、ことを特徴とする、付記 6 記載の文書収集方法。

(付記 1 0) 前記収集済み文書で用いられている参照表現に基づいて、前記収集済み文書群をまとめあげる、

ことを特徴とする付記 1 記載の文書収集方法。

【 0 1 6 7 】

(付記 1 1) 前記収集済み文書で用いられている参照表現に基づいて、前記収集済み文書にキーワードを付与する、

ことを特徴とする付記 1 記載の文書収集方法。

【 0 1 6 8 】

(付記 1 2) 前記参照表現が参照先文書に関係なく使用される参照表現の場合、キーワードとしない、

ことを特徴とする付記 1 1 記載の文書収集方法。

【 0 1 6 9 】

(付記 1 3) 前記参照表現が参照する相異なる文書数を計数し、前記相異なる文書数がある数以上である場合、前記参照表現をキーワードとしない、

ことを特徴とする付記 1 1 記載の文書収集方法。

【 0 1 7 0 】

(付記 1 4) 前記相異なる文書数がある数未満である場合、各収集済み文書で前記参照表現が参照されている回数である参照回数を計数し、

前記相異なる文書数及び前記参照回数に基づいて、前記参照表現をキーワードとするか否か判定する、

ことを特徴とする付記 1 1 記載の文書集収集方法。

【0 1 7 1】

(付記 1 5) 前記参照表現に基づくキーワードに、前記収集済み文書の本文から抽出したキーワード及び前記収集済み文書の URL から抽出したキーワードを組み合わせる、

ことを特徴とする付記 1 1 記載の文書集収集方法。

【0 1 7 2】

(付記 1 6) ネットワーク上のコミュニティに属する端末から文書を検索する検索方法であって、

文書を検索するための情報をサーバに送信し、

前記検索するための情報に基づいて前記コミュニティ内外に分けて検索した文書を、前記コミュニティにとっての有用度とともに受信する、

ことを特徴とする検索方法。

【0 1 7 3】

(付記 1 7) ネットワークから文書を収集する文書収集装置であって、

前記文書の参照関係に基づいて、次に収集すべき文書の候補である次収集候補を決定する次候補判定手段と、

前記文書のネットワーク上の場所を示す情報に基づいて前記文書が前記ネットワーク上のコミュニティ内の文書であるか否か判別するコミュニティ判別手段と

前記ネットワークから前記次収集候補を収集する文書収集手段と、を備え、  
前記文書収集手段は、前記コミュニティ内から所定数以上文書を収集した後、  
前記コミュニティ内外から文書を収集する、

ことを特徴とする文書収集装置。

【0 1 7 4】

(付記 1 8) ネットワークから文書を収集する文書収集装置であって、

ある分野に関する文書群である正例文書群及び前記分野と関連が少ない分野に

関する文書群である負例文書群の参照関係に基づいて、次に収集すべき文書の候補である次収集候補を決定する次候補判定手段と、

前記ネットワークから前記次収集候補を収集する文書収集手段とを備える、  
ことを特徴とする文書収集装置。

【0175】

(付記19) コンピュータに実行させることによって、ネットワークから文書を収集する制御を該コンピュータに行なわせるプログラムを記録した、コンピュータで読み取り可能な記録媒体であって、

前記文書の参照関係に基づいて、前記ネットワーク上のコミュニティ内から文書を所定数以上収集し、前記コミュニティから前記第1の所定数以上の文書を収集した後、収集済み文書の参照関係に基づいて、前記コミュニティ内外から文書を収集する、

ことを含む制御をコンピュータに行なわせるプログラムを記録した記録媒体。

【0176】

(付記20) コンピュータに実行させることによって、ネットワークからコンピュータに実行させることによって、ネットワークから文書を収集する制御を該コンピュータに行なわせるプログラムを記録した、コンピュータで読み取り可能な記録媒体であって、

ある分野に関する文書群である正例文書群及び前記分野と関連が少ない分野に関する文書群である負例文書群の参照関係に基づいて、前記分野に関する収集すべき文書を決定し、

前記ネットワークから前記収集すべき文書を収集する、

ことを含む制御をコンピュータに行なわせるプログラムを記録した記録媒体。

【0177】

(付記21) 搬送波に具現化された、ネットワークから文書を収集する制御をコンピュータに行わせるプログラムを表現するコンピュータ・データ・シグナルであって、前記プログラムは以下をコンピュータに実行させる、

前記文書の参照関係に基づいて、前記ネットワーク上のコミュニティ内から文書を所定数以上収集し、

前記コミュニティから前記所定数以上の文書を収集した後、収集済み文書の参照関係に基づいて、前記コミュニティ内外から文書を収集する、

【0178】

【発明の効果】

以上詳細に説明したように、本発明は、ある用途向けの文書を収集する際に、文書間の参照関係に基づいて収集すべき文書を決定し、決定された文書を収集することにより、言語に依存すること無く、迅速に用途にあった文書を選択して収集することが可能となる。

【0179】

また、参照表現に基づいて、収集済み文書をまとめあげ、各収集済み文書にキーワードを付与することにより、収集済み文書へのアクセスを容易とすることが可能となる。また、文書本文の内容を解析しないため、言語に依存せず、迅速にキーワードを付与することが可能となる。

【図面の簡単な説明】

【図1】

本発明の原理図である。

【図2】

第1実施形態に係わる文書収集装置の構成図である。

【図3】

URLテーブルのデータ構造の一例を示す図である。

【図4】

参照関係テーブルのデータ構造の一例を示す図である。

【図5】

参照表現テーブルのデータ構造の一例を示す図である。

【図6】

参照回数テーブルのデータ構造の一例を示す図である。

【図7】

第1実施形態に係わる文書収集装置が行う処理の大まかな流れを示すフローチャートである。

【図 8】

コミュニティ内の文書を収集する際に次収集候補を判定する処理を示すフローチャートである。

【図 9】

収集済み文書及び参照先文書をランキングする処理を示すフローチャートである。

【図 1 0】

収集済み文書を選別する処理を示すフローチャートである。

【図 1 1】

キーワード付与処理を示すフローチャートである。

【図 1 2】

収集した文書を提供する画面の一例を示す図である。

【図 1 3】

第 2 実施形態に係わる文書収集装置の構成図である。

【図 1 4】

$LT(S)$ 、 $LT(p)$ 、 $LS(d, X)$ 、 $LS(A, X)$  が意味する文書の参照関係を示す図である。

【図 1 5】

$CC(d, A, X)$  が意味する文書の参照関係を示す図である。

【図 1 6】

第 2 実施形態に係わる文書収集装置が行う処理を示すフローチャートである。

【図 1 7】

参照度を算出する式に含まれる各集合が意味する参照関係を示す図である。

【図 1 8】

共参照度を算出する式に含まれる各集合が意味する参照関係を示す図である。

【図 1 9】

第 2 実施形態の変形例に係わる文書収集装置が行う処理を示すフローチャートである。

【図 2 0】

文書収集装置の収集精度の実験結果を示す図（その 1）である。

【図 2 1】

文書収集装置の収集精度の実験結果を示す図（その 2）である。

【図 2 2】

情報処理装置の構成図である。

【図 2 3】

情報処理装置にプログラムやデータを供給する記録媒体、伝送信号及び伝送媒体を説明する図である。

【符号の説明】

1、1 0 0、2 0 0 文書収集装置

2 文書収集手段

3 参照関係抽出手段

4 コミュニティ判別手段

5 次候補判定手段

6 ランキング手段

7 U R L 判定手段

8 参照度／共参照度算出手段

9 まとめあげ手段

1 0 キーワード付与手段

2 0 収集済み文書群

2 1 次収集候補

2 2 文書間参照関係

2 3 収集文書ファイル

1 0 1 文書収集部

1 0 2 参照関係抽出部

1 0 3 コミュニティ判別部

1 0 4 候補判定部

1 0 5 ランキング部

1 0 6 まとめあげ部

- 107 キーワード付与部
- 120 URLテーブル
- 121 参照関係テーブル
- 122 参照表現テーブル
- 123 参照回数テーブル
- 130 優良コンテンツ
- 140 検索エンジン
- 141 索引
- 150 分類エンジン
- 160 サーバ
- 170 ブラウザ
- 180、181、182 画面
  - 201 参照度／共参照度テーブル
  - 210 分野別優良コンテンツ
  - 300 情報処理装置
    - 301 CPU
    - 302 メモリ
    - 303 入力装置
    - 304 出力装置
    - 305 外部記憶装置
    - 306 媒体駆動装置
    - 307 ネットワーク接続装置
    - 308 バス
    - 309 可搬記録媒体
    - 310 プログラム（データ）提供者
    - 311 回線

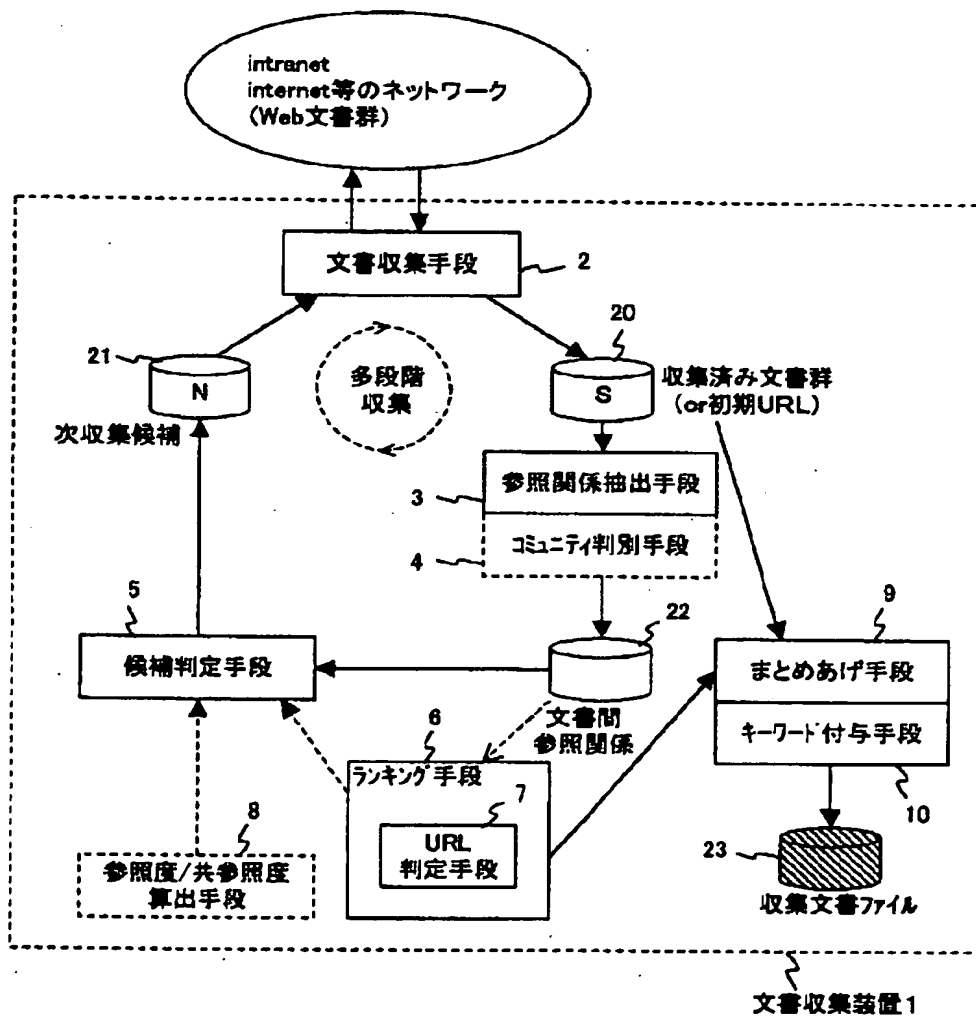
特 2 0 0 0 - 3 9 7 9 6 6



【書類名】 図面

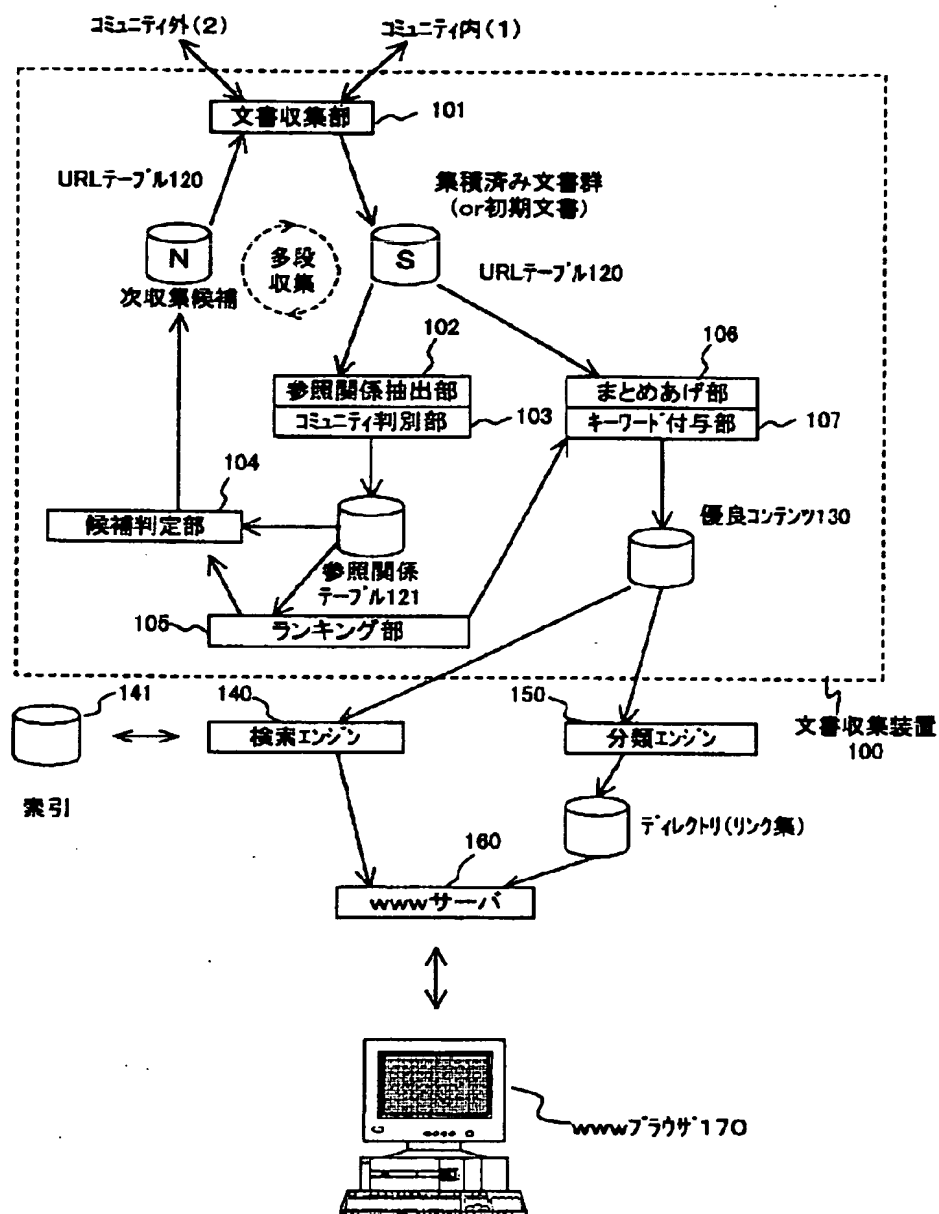
【図 1】

# 本発明の原理図



【図 2】

## 第1実施形態に係る文書収集装置の構成図



【図3】

## URLテーブルのデータ構造の一例を示す図

URLテーブル120

文書ID	URL	集積済みフラグ	コミュニティフラグ	重要度
doc1	http://www.abc.com/	1	1	....
doc2	http://www.klm.com/	1	0	....
....	....	....	....	....

【図 4】

参照関係テーブルのデータ構造の一例を示す図

参照関係テーブル121

参照元文書ID	参照先文書ID <sub>1</sub> (コミュニティ内)	参照先文書ID <sub>2</sub> (コミュニティ外)
doc1	doc2, doc3, doc4	doc5, doc6, doc7
doc2	doc4, doc8	doc9, doc10
....	....	....

【図 5】

参照表現テーブルのデータ構造の一例を示す図

参照表現テーブル122			
表現ID(w)	表現文字列	頻度(DF(w))	要否フラグ
rw1	ABC社	24	1
rw2	トップに戻る	28305	0
....	....	....	....

【図 6】

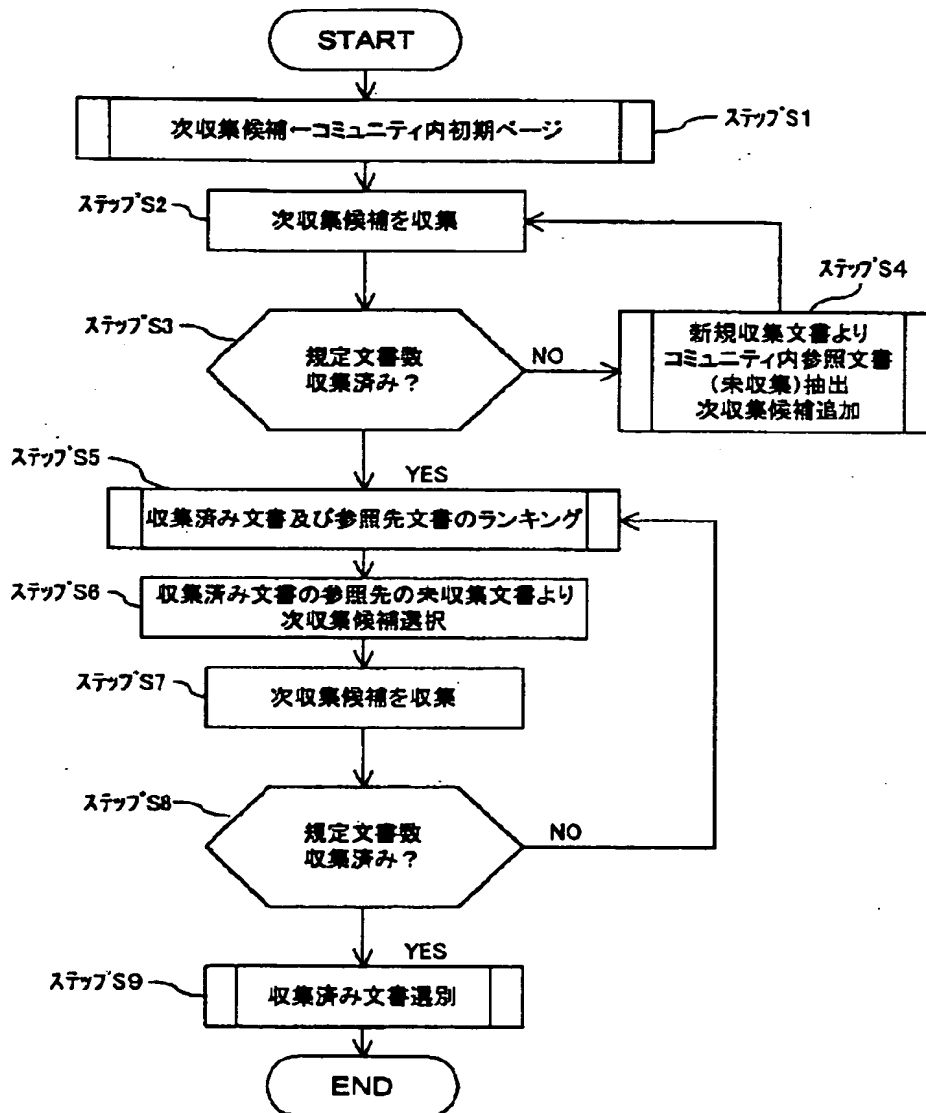
参照回数テーブルのデータ構造の一例を示す図

参照回数テーブル123

文書ID (d) \ 表現ID (w)	rw1	rw2	...	rwj	...	rw <sub>m</sub>
doc1	19	0	...		...	
doc1	0	2	...		...	...
doc1						
doc1				TF(i,j)		
doc1						
doc1	...	...	...		...	...

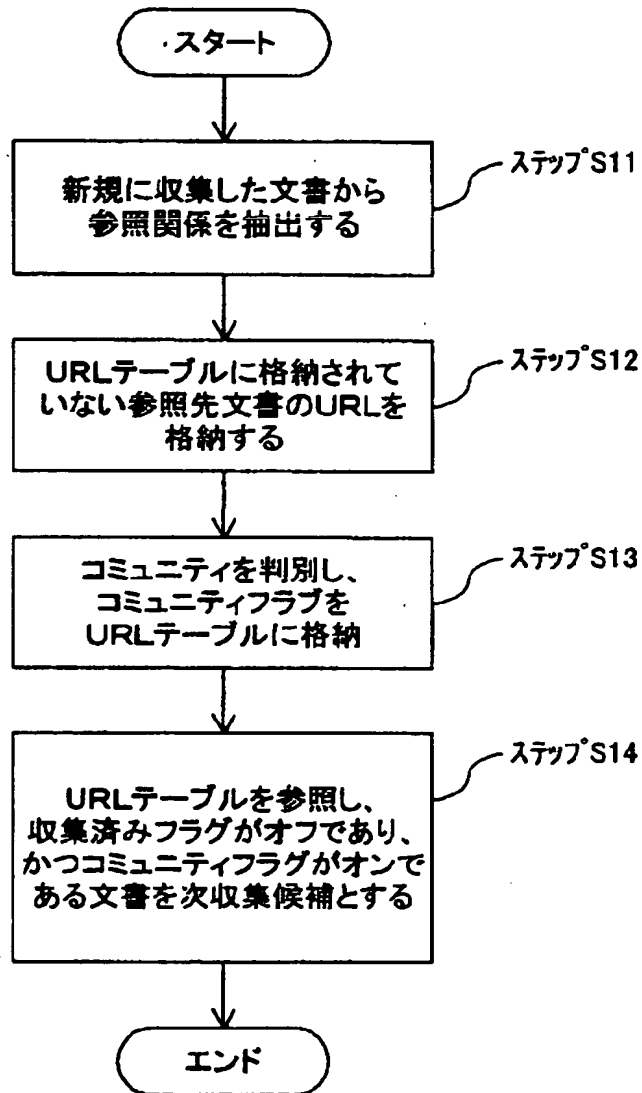
【図 7】

第一実施形態に係る文書収集装置が行う処理の  
大まかな流れを示すフローチャート



【図 8】

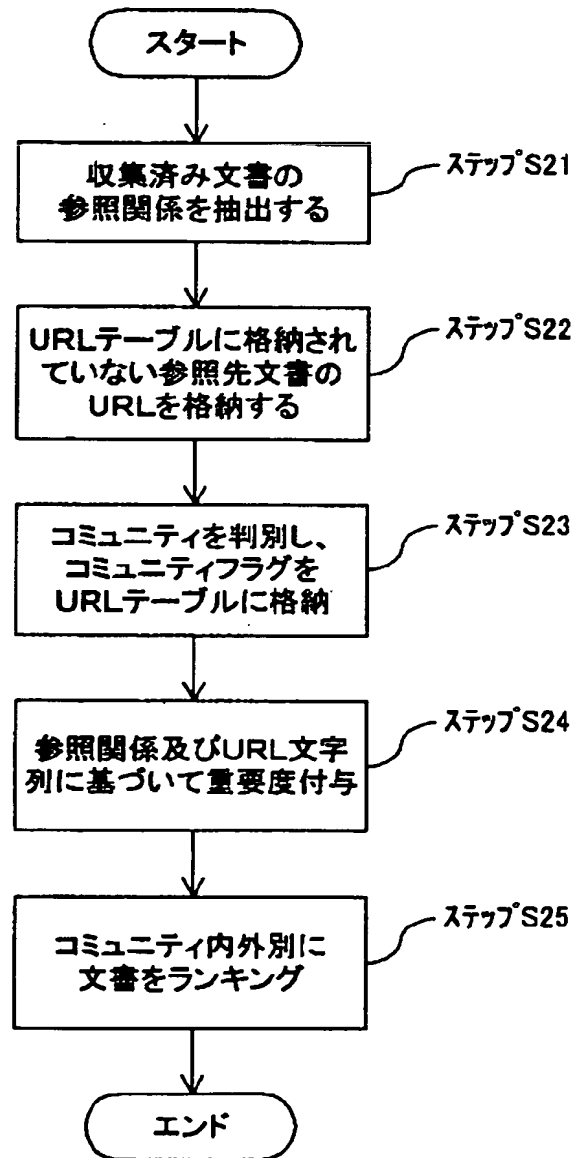
コミュニティ内の文書を収集する際に  
次収集候補を判定する処理を示すフローチャート





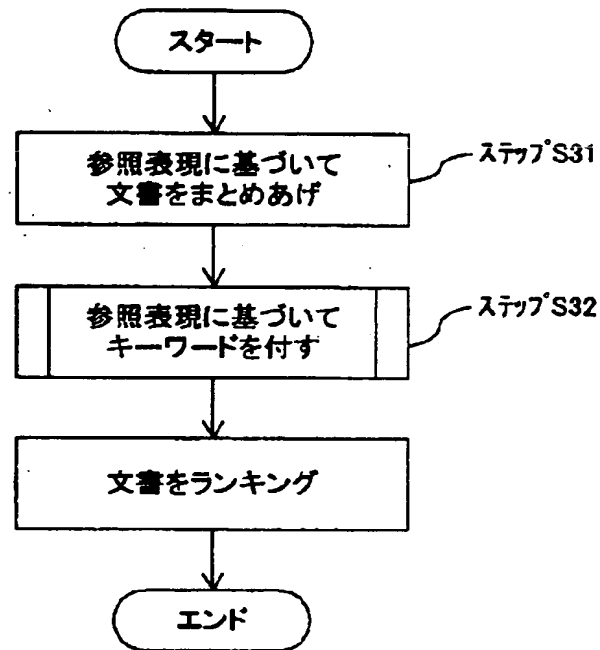
【図 9】

収集済み文書及び参照先文書を  
ランキングする処理を示すフローチャート



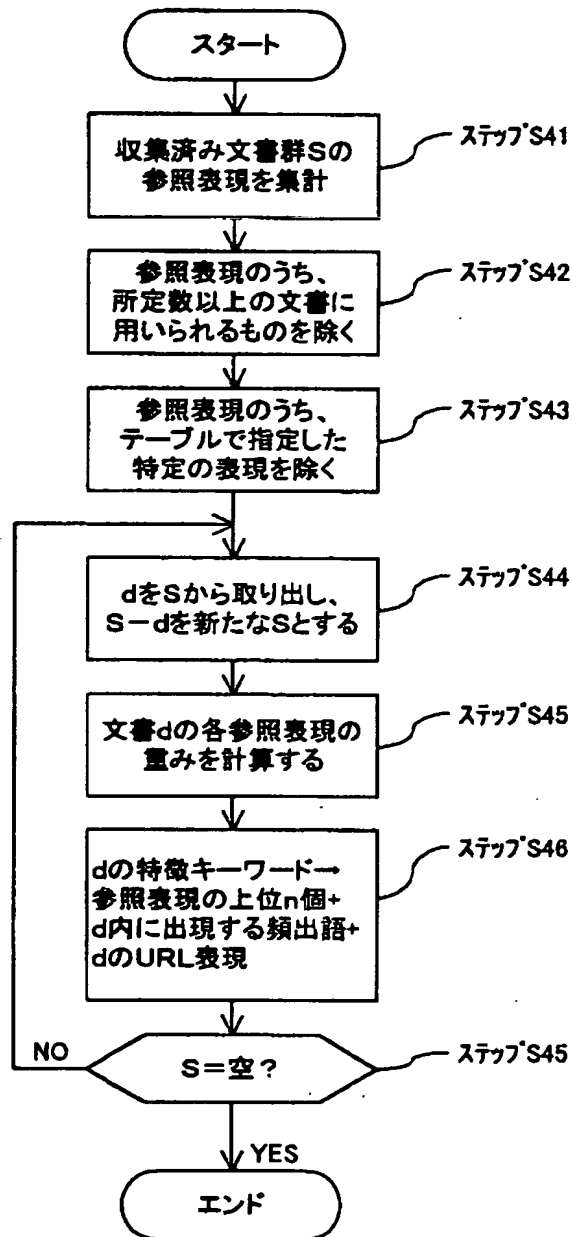
【図10】

収集済み文書を選別する処理を示すフローチャート



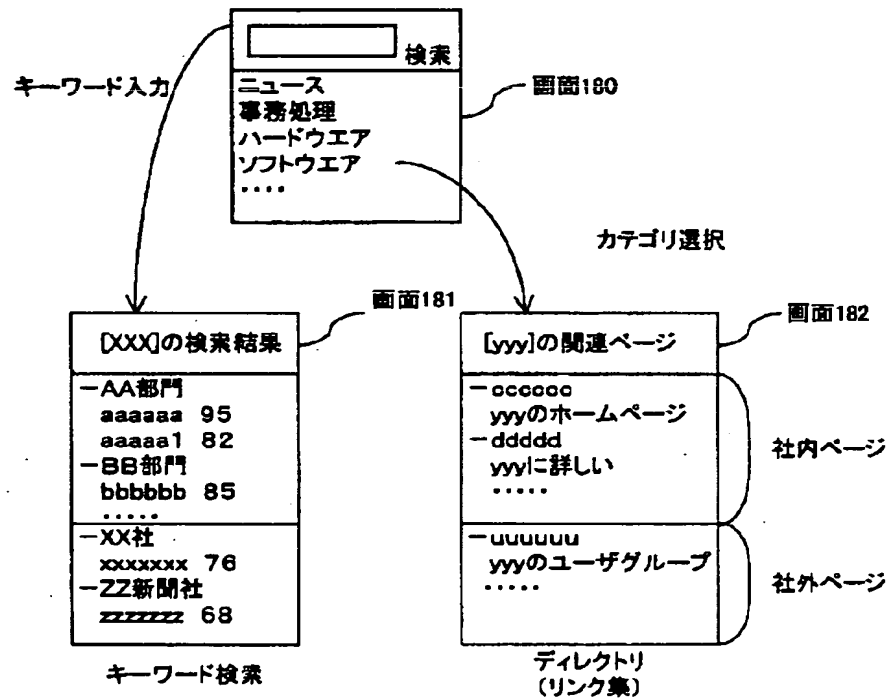
【図 11】

## キーワード付与処理を示すフローチャート



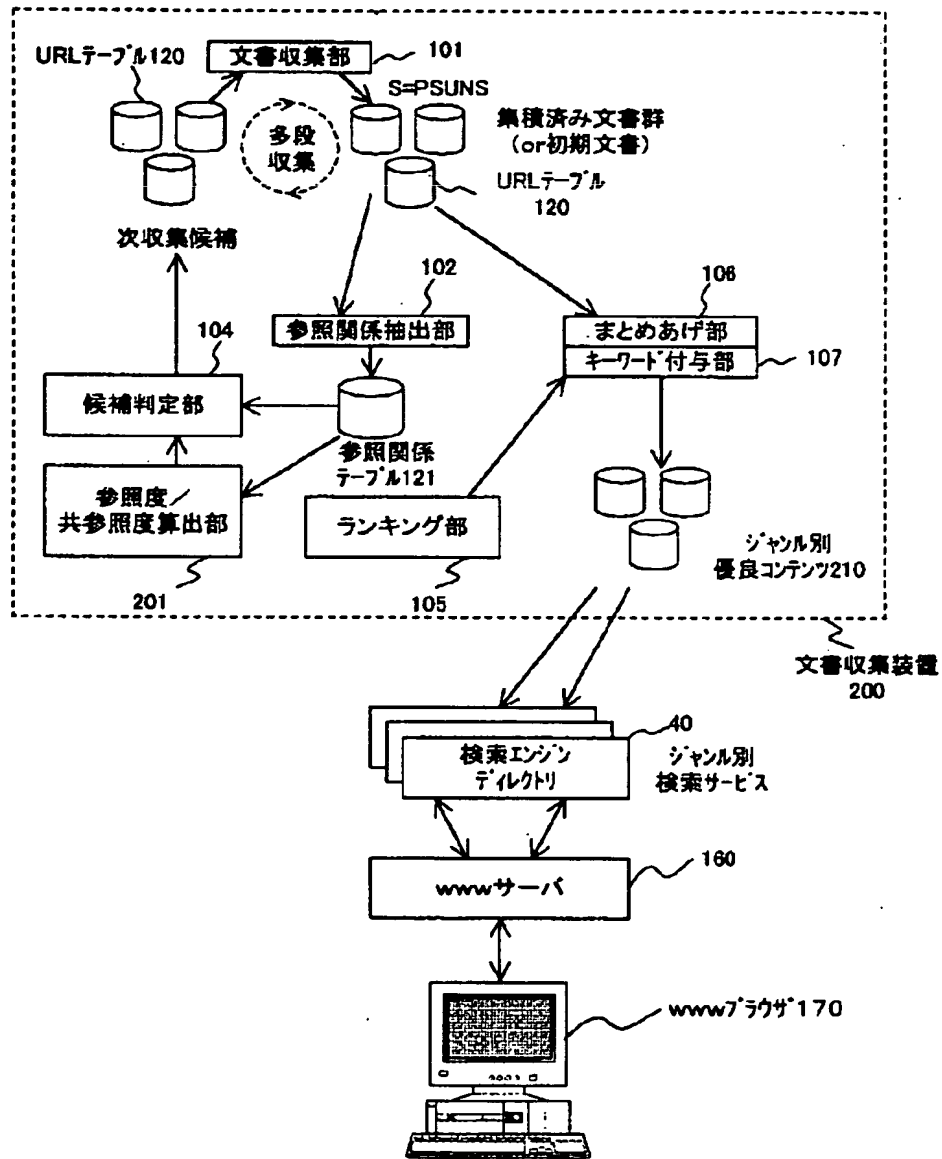
【図 12】

収集した文書を提供する画面の一例を示す図



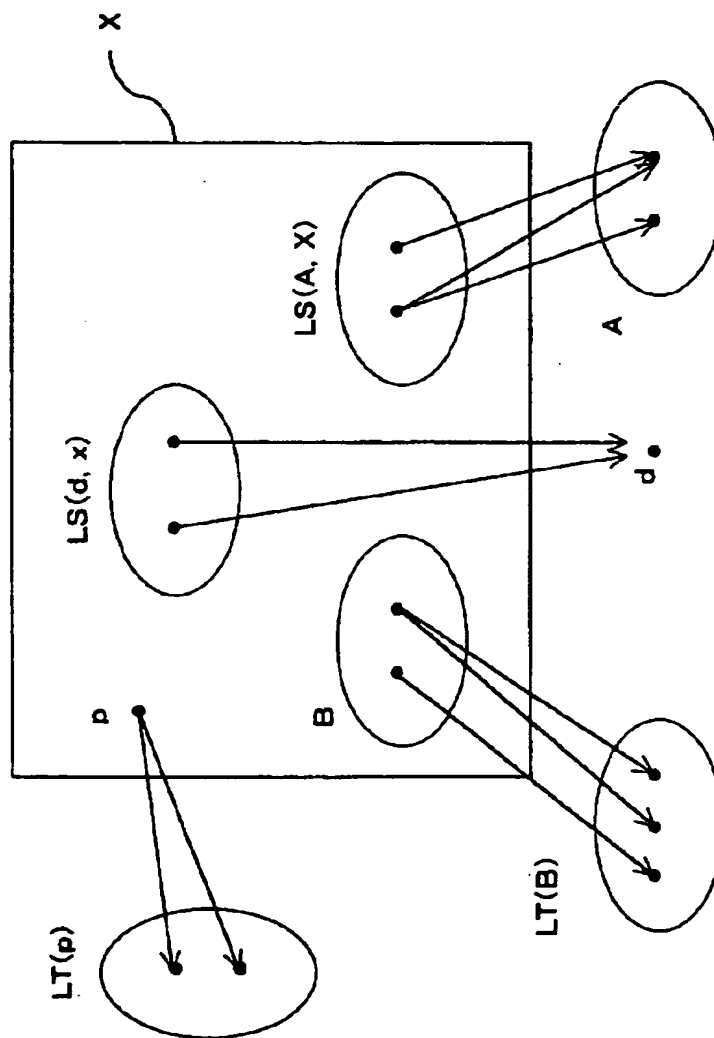
【図 13】

第2実施形態に係る文書収集装置の構成図



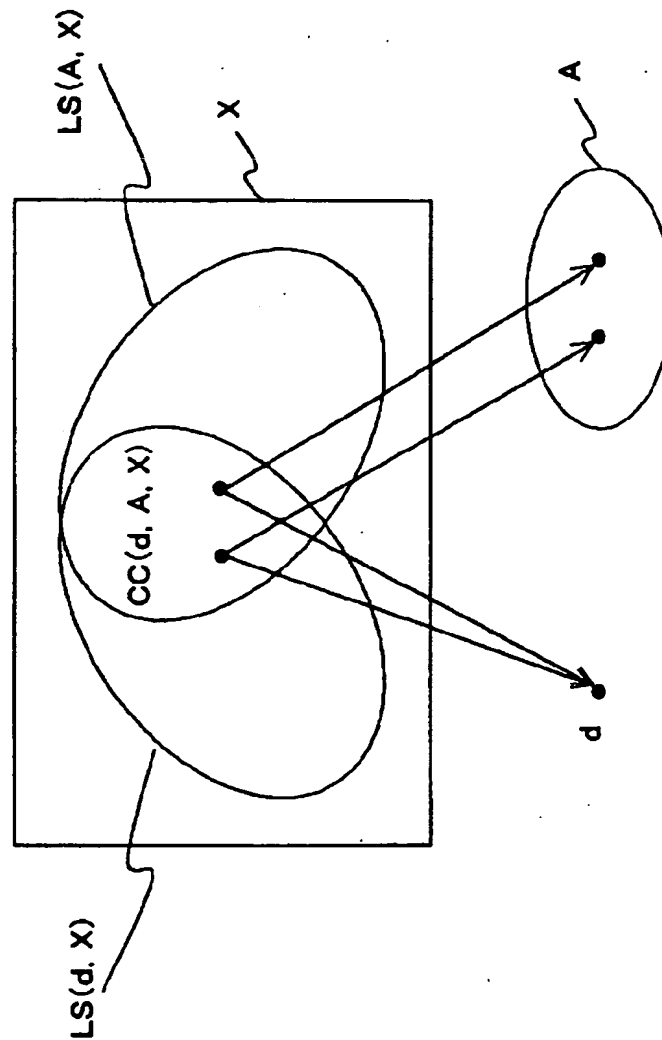
【図 14】

LT(s), LT(p), LS(d, x) 及び LS(A, X) が意味する、  
文書の参照関係を示す図



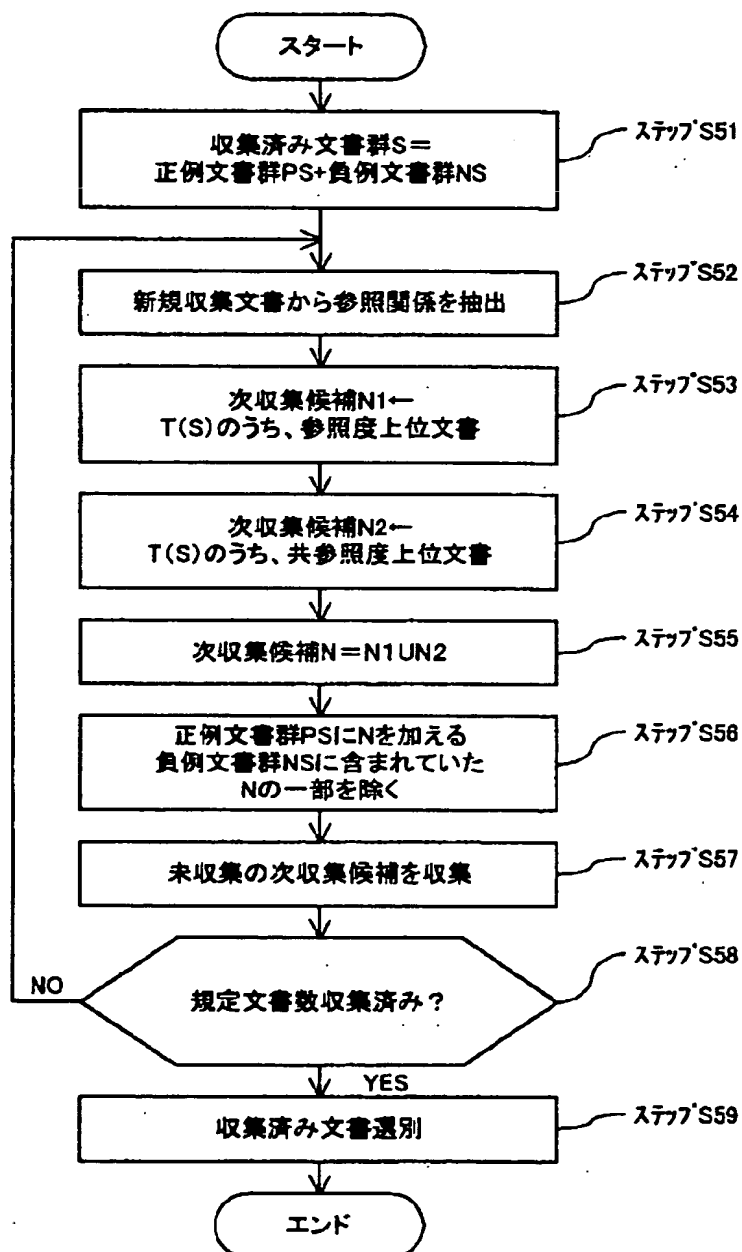
【図 15】

$CC(d, A, X)$  が意味する文書の参照関係を示す図



【図 16】

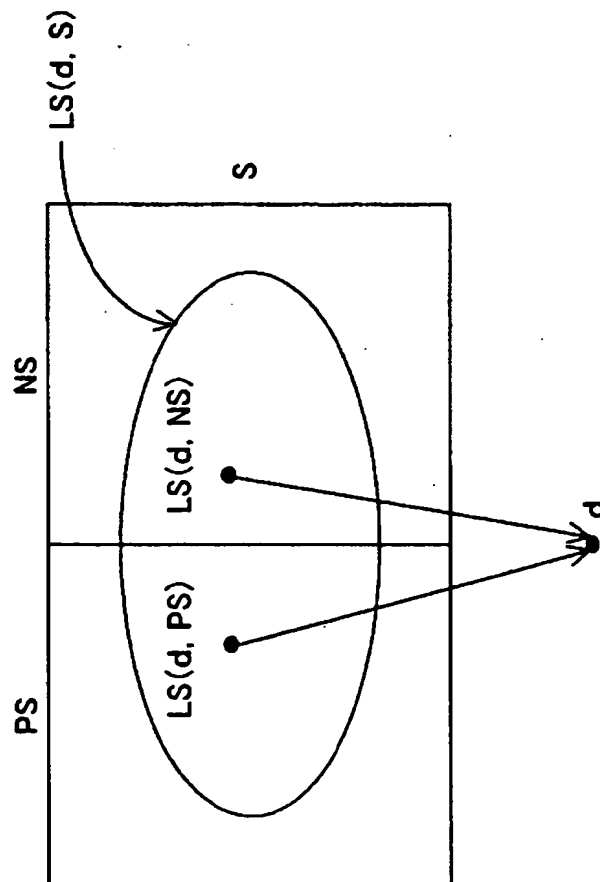
第2実施形態に係る文書収集装置が  
行う処理を示すフローチャート





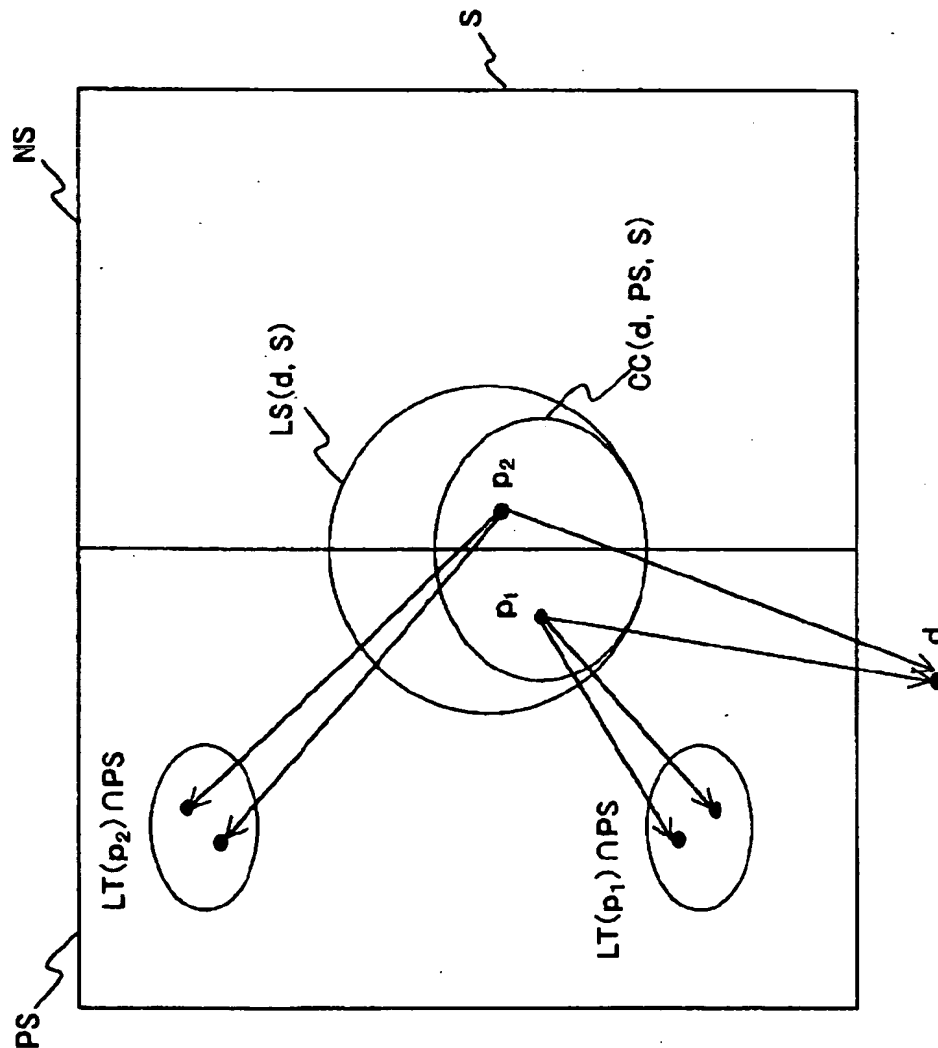
【図 17】

参照度を算出する式に含まれる  
各集合が意味する参照関係を示す図



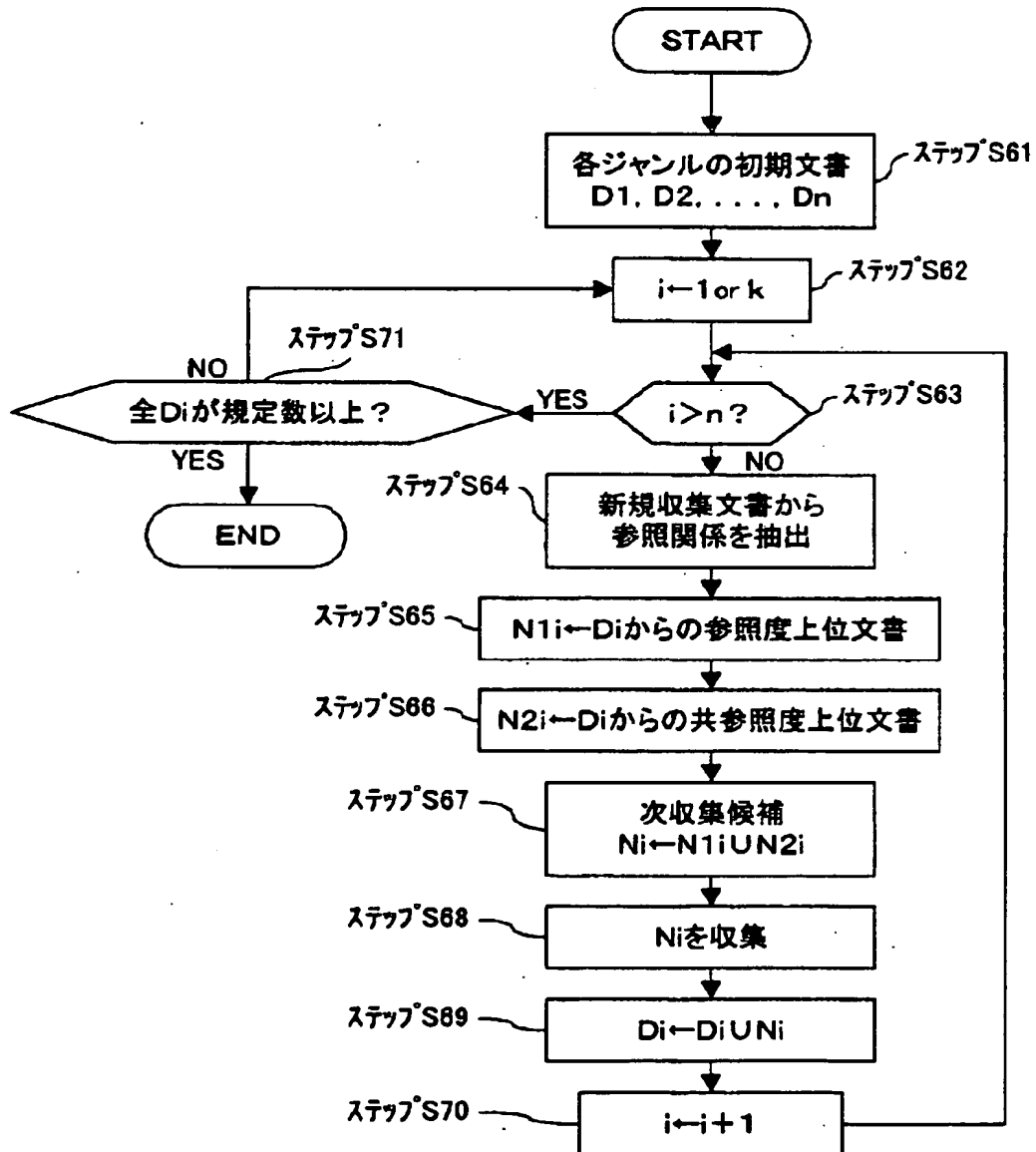
【図 18】

共参照度を算出する式に含まれる  
各集合が意味する参照関係を示す図



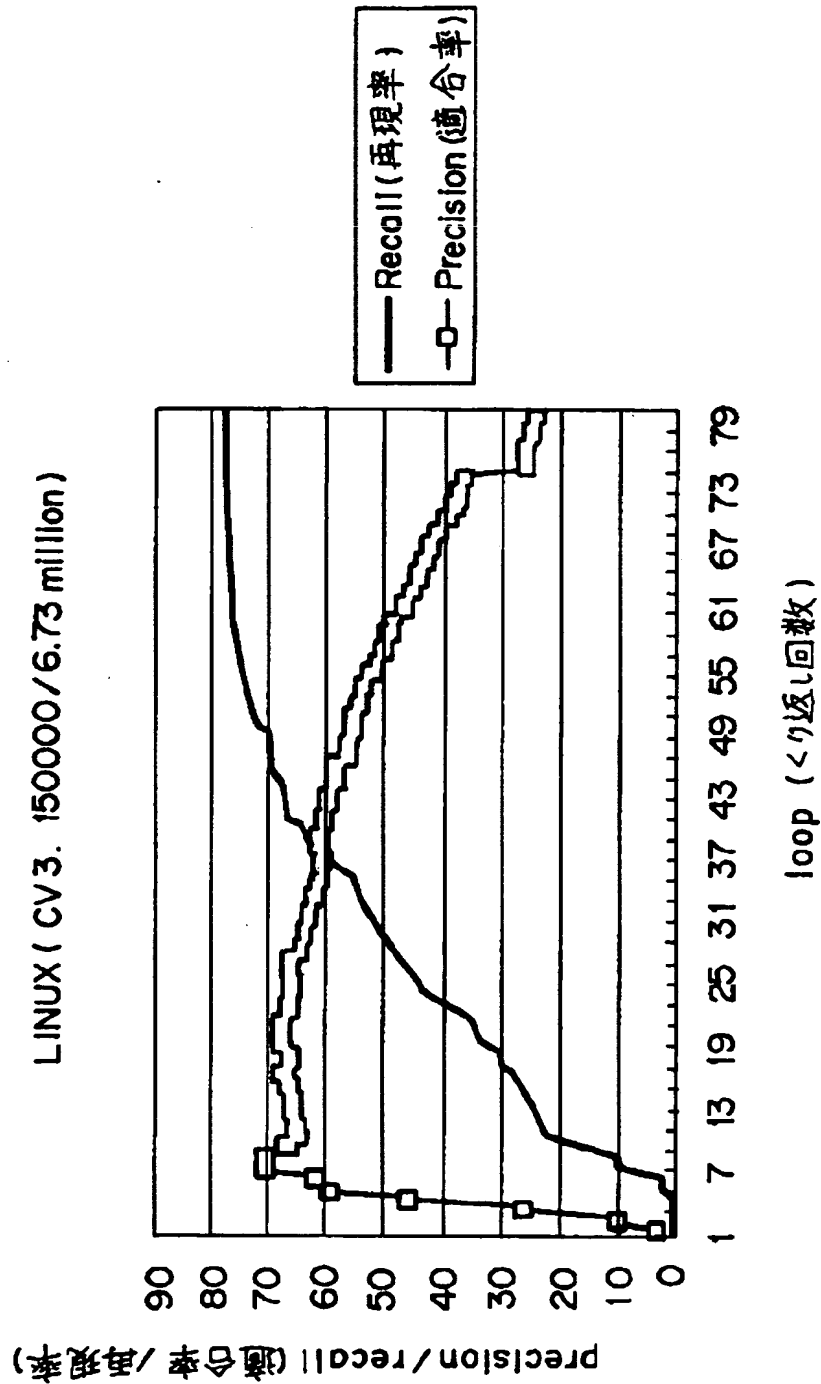
【図19】

第2実施形態の変形例に係る  
文書収集装置が行う処理を示すフローチャート



【図 20】

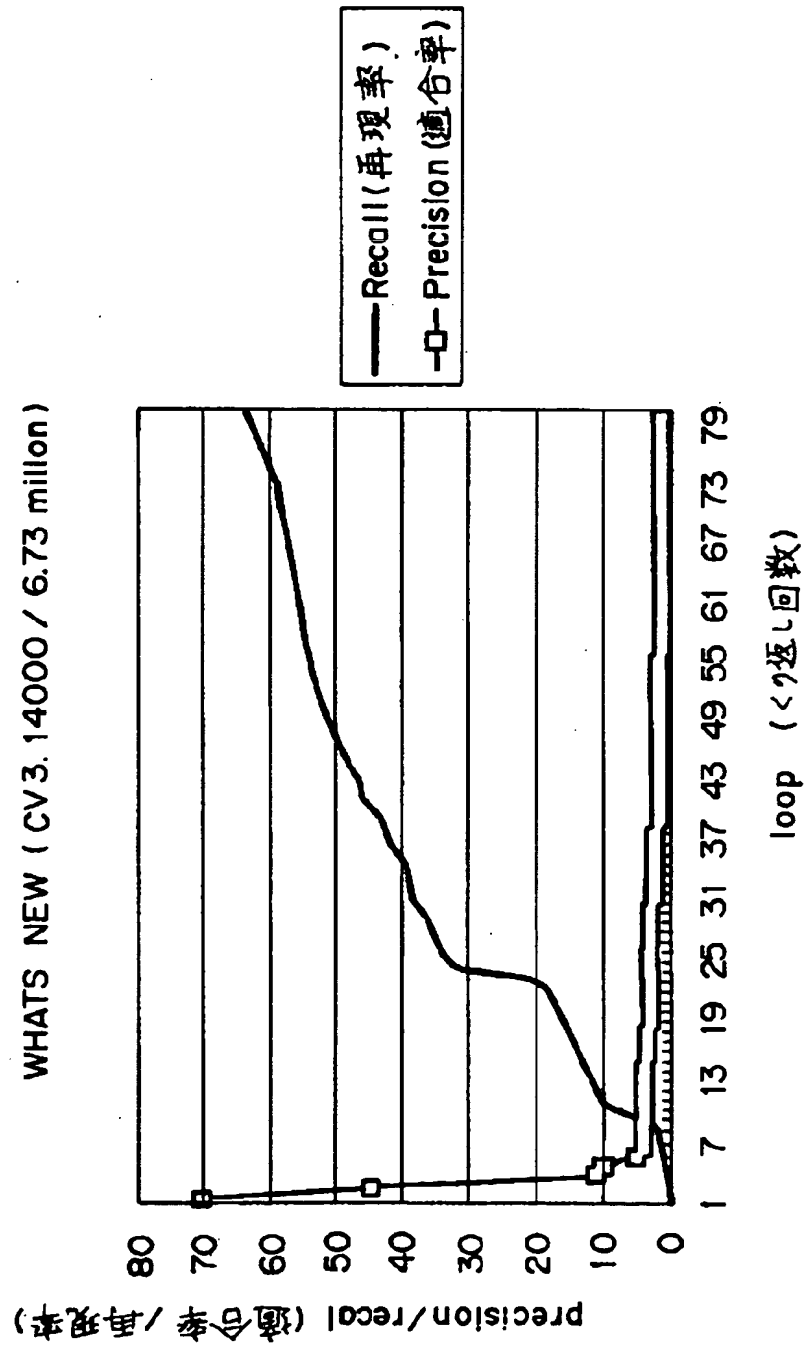
文書収集装置の収集精度の実験結果を示す図  
(その 1)



【図 2 1】

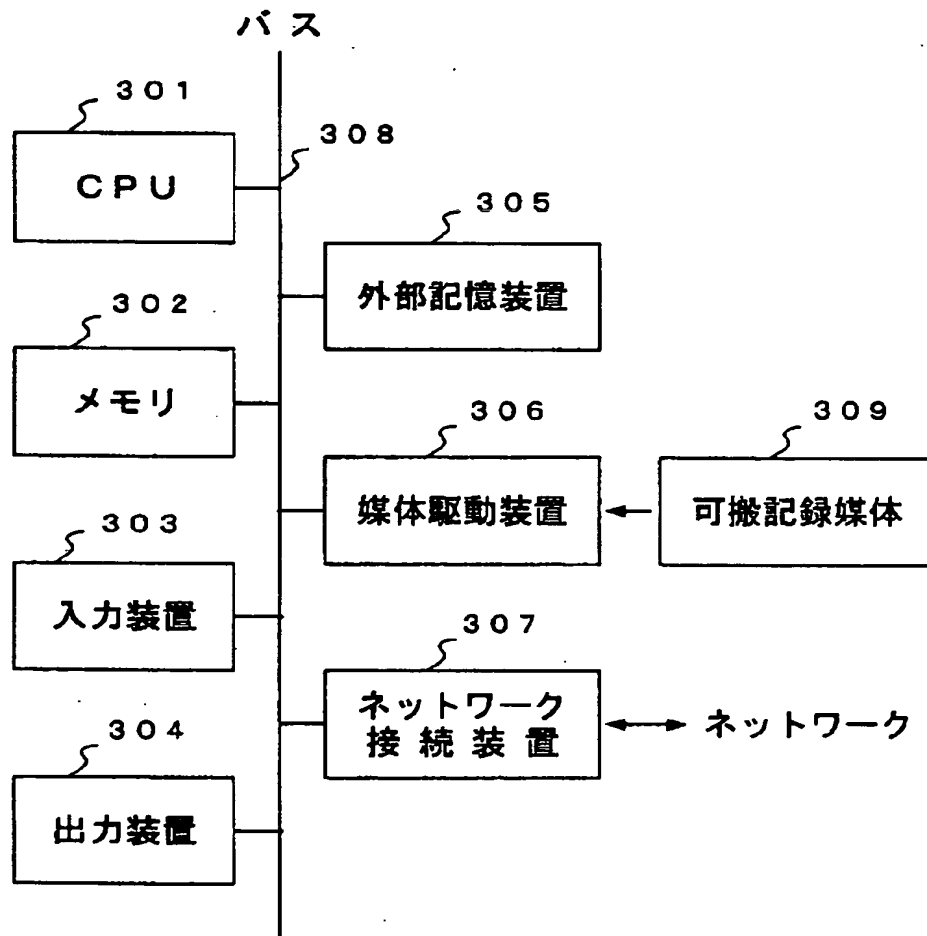
文書収集装置の収集精度の実験結果を示す図

(その 2)



【図 22】

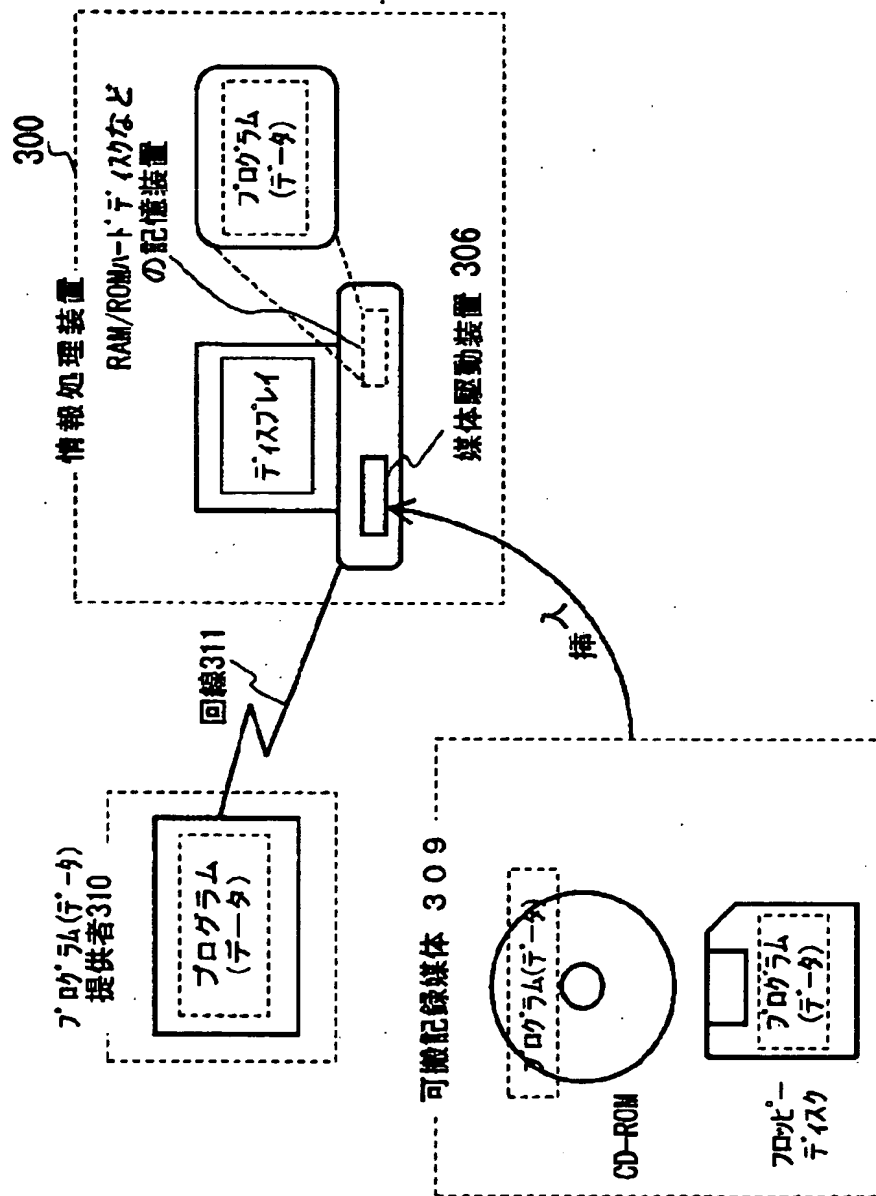
情 報 処 理 装 置 の 構 成 図



情報処理装置 300

【図 23】

情報処理装置にプログラムやデータを供給する記憶媒体、  
伝送信号、及び伝送媒体を説明する図



【書類名】 要約書

【要約】

【課題】 ネットワークから特定用途向けの文書を迅速に収集することを可能とする。

【解決手段】 ネットワークからコミュニティ向けの文書を収集する文書収集装置において、収集開始に先立って、収集の開始点となる初期文書群20を与える。参照関係抽出手段3は、初期文書群20から参照関係を抽出する。次候補判定手段5は、収集済みの文書群20の参照関係に基づいて、一定の条件を満たす未収集文書を次に収集すべき文書の候補である次収集候補21を判定する。文書収集手段2は、次収集候補21として判定された文書を収集し、収集済み文書20に加える。収集済み文書群20の文書数が規定された数以上でなければ、参照関係抽出手段3は、新規に収集された文書から更に参照関係を抽出し、上述の処理を繰り返す。このように、収集済み文書群20の文書数が規定された数以上になるまで文書の収集を繰り返す。

上記装置を、ネットワーク上のコミュニティにとって有用度の高い文書を収集するコミュニティ向けの文書収集装置として構成する場合、先に、文書収集手段2がネットワーク上のコミュニティ内から文書をまんべんなく収集した後、次候補判定手段5は、収集済み文書群20の参照関係に基づいてコミュニティ内外の文書から次収集候補21を決定する。

【選択図】 図1



職権訂正履歴（職権訂正）

特許出願の番号	特願2000-397966
受付番号	50001692245
書類名	特許願
担当官	濱谷 よし子 1614
作成日	平成13年 1月 4日

<訂正内容1>

訂正ドキュメント

明細書

訂正原因

職権による訂正

訂正メモ

【図面の簡単な説明】 【図15】を改行しました。

訂正前内容

【図14】

LT(S)、LT(p)、LS(d, X)、LS(A, X)が意味する文書の  
参照関係を示す図である。 【図15】

CC(d, A, X)が意味する文書の参照関係を示す図である。

【図16】

訂正後内容

【図14】

LT(S)、LT(p)、LS(d, X)、LS(A, X)が意味する文書の  
参照関係を示す図である。

【図15】

CC(d, A, X)が意味する文書の参照関係を示す図である。

【図16】

出 願 人 履 歴 情 報

識別番号 [000005223]

1. 変更年月日	1996年 3月26日
[変更理由]	住所変更
住 所	神奈川県川崎市中原区上小田中4丁目1番1号
氏 名	富士通株式会社